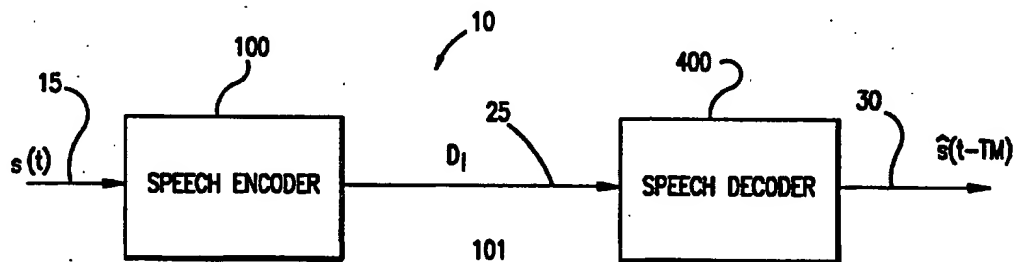


**PCT**WORLD INTELLECTUAL PROPERTY ORGANIZATION  
International Bureau

## INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification <sup>6</sup> : <b>G10L 3/02, 9/00</b>		A1	(11) International Publication Number: <b>WO 96/02050</b>
			(43) International Publication Date: 25 January 1996 (25.01.96)
(21) International Application Number: PCT/US95/08616			(81) Designated States: AM, AU, BB, BG, BR, BY, CA, CN, CZ, EE, FI, GE, HU, IS, JP, KG, KP, KR, KZ, LK, LR, LT, LV, MD, MG, MN, MX, NO, NZ, PL, RO, RU, SG, SI, SK, TJ, TM, TT, UA, UZ, VN, European patent (AT, BE, CH, DE, DK, ES, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, ML, MR, NE, SN, TD, TG), ARIPO patent (KE, MW, SD, SZ, UG).
(22) International Filing Date: 10 July 1995 (10.07.95)			
(30) Priority Data: 273,069 11 July 1994 (11.07.94) US			
(71) Applicant: VOXWARE, INC. [US/US]; 172 Tamarack Circle, Skillman, NJ 08558 (US).			
(72) Inventor: AGUILAR, Joseph, G.; 5148 West Wolfe Drive, Oak Lawn, IL 60453 (US).			
(74) Agents: MORRIS, Francis, E. et al.; Pennie & Edmonds, 1155 Avenue of the Americas, New York, NY 10036 (US).			Published With international search report.

(54) Title: HARMONIC ADAPTIVE SPEECH CODING METHOD AND SYSTEM



## (57) Abstract

A method and system is provided for encoding and decoding of speech signals at a low bit rate. The continuous input speech (15) is divided into voiced and unvoiced time segments of a predetermined length. The encoder of the system (100) uses a linear predictive coding model for the unvoiced speech segments and harmonic frequencies decomposition for the voiced speech segment. Only the harmonic frequencies are determined using the discrete fourier transform of the voiced speech segments. The decoder (400) synthesizes voice speech segments using the magnitudes of the transmitted harmonics and estimates the phase of each harmonic from the signal in the preceeding speech segments. Unvoiced speech segments are synthesized using linear prediction coding coefficients obtained from codebook entries for the poles of the LPC coefficient polynomial. Boundary conditions between voiced and unvoiced segments are established to insure amplitude and phase continuity for improved output speech quality.

AM

**FOR THE PURPOSES OF INFORMATION ONLY**

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AT	Austria	GB	United Kingdom	MR	Mauritania
AU	Australia	GE	Georgia	MW	Malawi
BB	Barbados	GN	Guinea	NE	Niger
BE	Belgium	GR	Greece	NL	Netherlands
BF	Burkina Faso	HU	Hungary	NO	Norway
BG	Bulgaria	IE	Ireland	NZ	New Zealand
BJ	Benin	IT	Italy	PL	Poland
BR	Brazil	JP	Japan	PT	Portugal
BY	Belarus	KE	Kenya	RO	Romania
CA	Canada	KG	Kyrgyzstan	RU	Russian Federation
CF	Central African Republic	KP	Democratic People's Republic of Korea	SD	Sudan
CG	Congo	KR	Republic of Korea	SE	Sweden
CH	Switzerland	KZ	Kazakhstan	SI	Slovenia
CI	Côte d'Ivoire	LI	Liechtenstein	SK	Slovakia
CM	Cameroon	LK	Sri Lanka	SN	Senegal
CN	China	LU	Luxembourg	TD	Chad
CS	Czechoslovakia	LV	Latvia	TG	Togo
CZ	Czech Republic	MC	Monaco	TJ	Tajikistan
DE	Germany	MD	Republic of Moldova	TT	Trinidad and Tobago
DK	Denmark	MG	Madagascar	UA	Ukraine
ES	Spain	ML	Mali	US	United States of America
FI	Finland	MN	Mongolia	UZ	Uzbekistan
FR	France			VN	Viet Nam
GA	Gabon				

- 1 -

HARMONIC ADAPTIVE SPEECH CODING METHOD AND SYSTEMBACKGROUND OF THE INVENTION

5       The present invention relates to speech processing and more specifically to a method and system for low bit rate digital encoding and decoding of speech using harmonic analysis and synthesis of the voiced portions and predictive coding of the unvoiced.  
10 portions of the speech.

Reducing the bit rate needed for storage and transmission of a speech signal while preserving its perceptual quality is among the primary objectives of  
15 modern digital speech processing systems. In order to meet these contradicting requirements various models of the speech formation process have been proposed in the past. Most frequently, speech is modeled on a short-time basis as the response of a linear system  
20 excited by a periodic impulse train for voiced sounds or random noise for the unvoiced sounds. For mathematical convenience, it is assumed that the speech signal is stationary within a given short time segment, so that the continuous speech is represented  
25 as an ordered set of distinct voiced and unvoiced speech segments.

Voiced speech segments, which correspond to vowels in a speech signal, typically contribute most  
30 to the intelligibility of the speech which is why it is important to accurately represent these segments. However, for a low-pitched voice, a set of more than 80 harmonic frequencies ("harmonics") may be measured within a voiced speech segment within a 4 kHz  
35 bandwidth. Clearly, encoding information about all harmonics of such segment is only possible if a large

- 2 -

number of bits is used. Therefore, in applications where it is important to keep the bit rate low, simplified speech models need to be employed.

5

One conventional solution for encoding speech at low bit rates is based on a sinusoidal speech representation model. U.S. Patent No. 5,054,072 to McAuley for example describes a method for speech  
10 coding which uses a pitch extraction algorithm to model the speech signal by means of a harmonic set of sinusoids that serve as a "perceptual" best fit to the measured sinusoids in a speech segment. The system generally attempts to encode the amplitude envelope of  
15 the speech signal by interpolating this envelope with a reduced set of harmonics. In a particular embodiment, one set of frequencies linearly spaced in the baseband (the low frequency band) and a second set of frequencies logarithmically spaced in the high  
20 frequency band are used to represent the actual speech signal by exploiting the correlation between adjacent sinusoids. A pitch adaptive amplitude coder is then used to encode the amplitudes of the estimated harmonics. The proposed method, however, does not  
25 provide accurate estimates, which results in distortions of the synthesized speech.

The McAuley patent also provides a model for predicting the phases of the high frequency harmonics  
30 from the set of coded phases of the baseband harmonics. The proposed phase model, however, requires a considerable computational effort and furthermore requires the transmission of additional bits to encode the baseband harmonics phases so that very low bit  
35 rates may not be achieved using the system.

- 3 -

U.S. Patent No. 4,771,465 describes a speech analyzer and synthesizer system using a sinusoidal encoding and decoding technique for voiced speech segments and noise excitation or multipulse excitation for unvoiced speech segments. In the process of encoding the voiced segments a fundamental subset of harmonic frequencies is determined by a speech analyzer and is used to derive the parameters of the remaining harmonic frequencies. The harmonic amplitudes are determined from linear predictive coding (LPC) coefficients. The method of synthesizing the harmonic spectral amplitudes from a set of LPC coefficients, however, requires extensive computations using high precision floating point arithmetic and yields relatively poor quality speech.

U.S. Patent Nos. 5,226,108 and 5,216,747 to Hardwick et al. describe an improved pitch estimation method providing sub-integer resolution. The quality of the output speech according to the proposed method is improved by increasing the accuracy of the decision as to whether given speech segment is voiced or unvoiced. This decision is made by comparing the energy of the current speech segment to the energy of the preceding segments. Furthermore, harmonic frequencies in voiced speech segments are generated using a hybrid approach in which some harmonics are generated in the time domain while the remaining harmonics are generated in the frequency domain. According to the proposed method, a relatively small number of low-frequency harmonics are generated in the time domain and the remaining harmonics are generated in the frequency domain. Voiced harmonics generated in the frequency domain are then frequency scaled, transformed into the time domain using a discrete

- 4 -

Fourier transform (DFT), linearly interpolated and finally time scaled. The proposed method generally does not allow accurate estimation of the amplitude and phase information for all harmonics and is computationally expensive.

U.S. Patent No. 5,226,084 also to Hardwick et al. describes methods for quantizing speech while preserving its perceptual quality. To this end, harmonic spectral amplitudes in adjacent speech segments are compared and only the amplitude changes are transmitted to encode the current frame. A segment of the speech signal is transformed to the frequency domain to generate a set of spectral amplitudes. Prediction spectral amplitudes are then computed using interpolation based on the actual spectral amplitudes of at least one previous speech segment. The differences between the actual spectral amplitudes for the current segment and the prediction spectral amplitudes derived from the previous speech segments define prediction residuals which are encoded. The method reduces the required bit rate by exploiting the amplitude correlation between the harmonic amplitudes in adjacent speech segments, but is computationally expensive.

While the prior art discloses some advances toward achieving a good quality speech at a low bit rate, it is perceived that there exists a need for improved methods for encoding and decoding of speech at such low bit rates. More specifically, there is a need to obtain accurate estimates of the amplitudes of the spectral harmonics in voiced speech segments in a computationally efficient way and to develop a method and system to synthesize such voiced speech segments

- 5 -

without the requirement to store or transmit separate phase information.

5

10

15

20

25

30

35

- 6 -

SUMMARY OF THE INVENTION

Accordingly, it is an object of the present invention to provide a low bit-rate method and system  
5 for encoding and decoding of speech signals using adaptive harmonic analysis and synthesis of the voiced portions and predictive coding of the unvoiced portions of the speech signal.

10 It is another object of the present invention to provide a super resolution harmonic amplitude estimator for approximating the speech signal in a voiced time segment as a set of harmonic frequencies.

15 It is another object of the present invention to provide a novel phase compensated harmonic synthesizer to synthesize speech in voiced segments from a set of harmonic amplitudes and combine the generated speech segment with adjacent voiced or unvoiced speech  
20 segments with minimized amplitude and phase distortions to obtain good quality speech at a low bit rate.

These and other objectives are achieved in  
25 accordance with the present invention by means of a novel encoder/decoder speech processing system in which the input speech signal is represented as a sequence of time segments (also referred to as frames), where the length of the time segments is  
30 selected so that the speech signal within each segment is relatively stationary. Thus, dependent on whether the signal in a time segment represents voiced (vowels) or unvoiced (consonants) portions of the speech, each segment can be classified as either being  
35 voiced or unvoiced.



- 7 -

In the system of the present invention the continuous input speech signal is digitized and then divided into segments of predetermined length. For each input segment a determination is next made as to whether it is voiced or unvoiced. Dependent on this determination, each time segment is represented in the encoder by a signal vector which contains different information. If the input segment is determined to be unvoiced, the actual speech signal is represented by the elements of a linear predictive coding vector. If the input segment is voiced, the signal is represented by the elements of a harmonic amplitudes vector. Additional control information including the energy of the segment and the fundamental frequency in voiced segments is attached to each predictive coding and harmonic amplitudes vector to form data packets. The ordered sequence of data packets completely represents the input speech signal. Thus, the encoder of the present invention outputs a sequence of data packets which is a low bit-rate digital representation of the input speech.

More specifically, after the analog input speech signal is digitized and divided into time segments, the system of the present invention determines whether the segment is voiced or unvoiced using a pitch detector to this end. This determination is made on the basis of the presence of a fundamental frequency in the speech segment which is detected by the pitch detector. If such fundamental frequency is detected, the pitch detector estimates its frequency and outputs a flag indicating that the speech segment is voiced.

If the segment is determined to be unvoiced, the system of the present invention computes the roots of

- 8 -

a characteristic polynomial with coefficients which are the LPC coefficients for the speech segment. The computed roots are then quantized and replaced by a  
5 quantized vector codebook entry which is representative of the unvoiced time segment. In a specific embodiment of the present invention the roots of the characteristic polynomial may be quantized using a neural network linear vector quantizer (LVQ1).

10

If the speech segment is determined to be voiced, it is passed to a novel super resolution harmonic amplitude estimator which estimates the amplitudes of the harmonic frequencies of the speech segment and  
15 outputs a vector of normalized harmonic amplitudes representative of the speech segment.

A parameter encoder next generates for each time segment of the speech signal a data packet, the  
20 elements of which contain information necessary to restore the original signal segment. For example, a data packet for an unvoiced speech segment comprises control information, a flag indicating that the segment is unvoiced, the total energy of the segment  
25 or the prediction error power, and the elements of the codebook entry defining the roots of the LPC coefficient polynomial. On the other hand, a data packet for a voiced speech segment comprises control information, a flag indicating that the segment is  
30 voiced, the sum total of the harmonic amplitudes of the segment, the fundamental frequency and a set of estimated normalized harmonic amplitudes. The ordered sequence of data packets at the output of the parameter encoder is ready for storage or transmission  
35 of the original speech signal.

- 9 -

At the synthesis side, a decoder receives the ordered sequence of data packets representing unvoiced and voiced speech signal segments. If the  
5   voiced/unvoiced flag indicates that a data packet represents an unvoiced time segment, the transmitted quantized pole vector is used as an index into a pole codebook to determine the LPC coefficients of the unvoiced synthesis (prediction) filter. A gain  
10   adjusted white noise generator is then used as the input of the synthesis filter to reconstruct the unvoiced speech segment.

If the data packet flag indicates that a segment  
15   is voiced, a novel phase compensated harmonic synthesizer is used to synthesize the voiced speech segment and provide amplitude and phase continuity to the signal of the preceding speech segment. Specifically, using the harmonic amplitudes vector of  
20   the voiced data packet, the phase compensated harmonic synthesizer computes the conditions required to insure amplitude and phase continuity between adjacent voiced segments and computes the parameters of the voiced to unvoiced or unvoiced to voiced speech segment  
25   transitions. The phases of the harmonic frequencies in a voiced segment are computed from a set of equations defining the phases of the harmonic frequencies in the previous segment. The amplitudes of the harmonic frequencies in a voiced segment are  
30   determined from a linear interpolation of the received amplitudes of the current and the previous time segments. Continuous boundary conditions between signal transitions at the ends of the segment are finally established before the synthesized signal is  
35   passed to a digital-to-analog converter to reproduce the original speech.

- 10 -

BRIEF DESCRIPTION OF THE DRAWINGS

The invention will be next be described in detail by reference to the following drawings in which:

5           Fig. 1 is a block diagram of the speech processing system of the present invention.

          Fig. 2 is a schematic block diagram of the encoder used in the system of Fig. 1.

10           Fig. 3 illustrates the signal sequences of the digitized input signal  $s(n)$  which define delayed speech vectors  $S_M(M)$  and  $S_{N-M}(N)$  used in the encoder of Fig. 2.

          Figs. 4 and 5 are schematic diagrams of the transmitted parameters in an unvoiced and in a voiced data packet, respectively.

15           Fig. 6 is a flow diagram of the super resolution harmonic amplitude estimator (SRHAE) used in the encoder in Fig. 2.

          Figs. 7A is a graph of the actual and the estimated harmonic amplitudes in a voiced speech segment.

          Fig. 7B illustrates the normalized estimation error in percent % dB for the harmonic amplitudes of the speech segment in Fig. 7A.

25           Fig. 8 is a schematic block diagram of the decoder used in the system of Fig. 1.

          Fig. 9 is a flow diagram of the phase compensated harmonic synthesizer in Fig. 8.

30           Figs. 10 A, B illustrate of the harmonics matching problem in the system of the present invention.

          Fig. 11 is a flow diagram of the voiced to voiced speech synthesis algorithm.

35           Fig. 12 is a flow diagram of the unvoiced to voiced speech synthesis algorithm.

- 11 -

Fig. 13 is a flow diagram of the initialization of the system with the parameters of the previous speech segment.

5

10

15

20

25

30

35

- 12 -

DETAILED DESCRIPTION OF THE INVENTION

During the course of the description like numbers will be used to identify like elements shown in the figures. Bold face letters represent vectors, while vector elements and scalar coefficients are shown in standard print.

Fig. 1 is a block diagram of the speech processing system 10 for encoding and decoding speech in accordance with the present invention. Analog input speech signal  $s(t)$ , 15 from an arbitrary voice source is received at encoder 100 for subsequent storage or transmission over a communications channel. Encoder 100 digitizes the analog input speech signal 15, divides the digitized speech sequence into speech segments and encodes each segment into a data packet 25 of length  $I$  information bits. The encoded speech data packets 25 are transmitted over communications channel 101 to decoder 400. Decoder 400 receives data packets 25 in their original order to synthesize a digital speech signal which is then passed to a digital-to-analog converter to produce a time delayed analog speech signal 30, denoted  $s(t-T_m)$ , as explained in detail next.

A. The Encoder Block

Fig. 2 illustrates the main elements of encoder 100 and their interconnections in greater detail. Blocks 105, 110 and 115 perform signal pre-processing to facilitate encoding of the input speech. In particular, analog input speech signal 15 is low pass filtered in block 105 to eliminate frequencies outside the human voice range. Low pass filter (LPF) 105 has a cutoff frequency of about 4 KHz which is adequate for the purpose. The low pass filtered analog signal

- 13 -

is then passed to analog-to-digital converter 110 where it is sampled and quantized to generate a digital signal  $s(n)$  suitable for subsequent processing. Analog-to-digital converter 110 preferably operates at a sampling frequency  $f_s = 8$  KHz which, in accordance with the Nyquist criterion, corresponds to twice the highest frequency in the low pass filtered analog signal  $s(t)$ . It will be appreciated that other sampling frequencies may be used as long as they satisfy the Nyquist criterion. Finally, digital input speech signal  $s(n)$  is passed through a high pass filter (HPF) 115 which has a cutoff frequency of about 100 Hz in order to eliminate any low frequency noise, such as 60 Hz AC voltage interference.

The filtered digital speech signal  $s(n)$  is next divided into time segments of a predetermined length in frame segmenters 120 and 125. Digital speech signal  $s(n)$  is first buffered in frame segmenter 120 which outputs a delayed speech vector  $S_M(M)$  of length  $M$  samples. Frame segmenter 120 introduces a time delay of  $M$  samples between the current sample of speech signal  $s(n)$  and the output speech vector  $S_M(M)$ . In a specific embodiment of the present invention, the length  $M$  is selected to be about 160 samples which corresponds to 20 msec of speech at a 8 KHz sampling frequency. This length of the speech segment has been determined to present a good compromise between the requirement to use relatively short segments as to keep the speech signal roughly stationary, and the efficiency of the coding system which generally increases as the delay becomes greater. Dependent on the desired temporal resolution, the delay between

- 14 -

time segments can be set to other values, such as 50, 100 or 150 samples.

5       A second frame segmenter 125 buffers  $N-M$  samples into a vector  $S_{N-M}(N)$ , the last element of which is delayed by  $N$  samples from the current speech sample  $s(n)$ . Fig. 3 illustrates the relationship between delayed speech vectors  $S_M(M)$ ,  $S_{N-M}(N)$  and the digital  
10 input speech signal  $s(n)$ . The function of the delayed vector  $S_{N-M}(N)$  will be described in more detail later.

The step following the segmentation of digital input signal  $s(n)$  is to decide whether the current  
15 segment is voiced or unvoiced, which decision determines the type of applied signal processing. Speech is generally classified as voiced if a fundamental frequency is imported to the air stream by the vocal cords of the speaker. In such case the  
20 speech signal is modeled as a superposition of sinusoids which are harmonically related to the fundamental frequency as discussed in more detail next. The determination as to whether a speech segment is voiced or unvoiced, and the estimation of  
25 the fundamental frequency can be obtained in a variety of ways known in the art as pitch detection algorithms.

In the system of the present invention, pitch  
30 detection block 155 determines whether the speech segment associated with delayed speech vector  $S_M(M)$  is voiced or unvoiced. In a specific embodiment, block 155 employs the pitch detection algorithm described in Y. Medan et al., "Super Resolution Pitch Determination  
35 of Speech Signals", IEEE Trans. on Signal Processing, Vol. 39, pp 40-48, June 1991, which is incorporated



- 15 -

herein by reference. It will be appreciated that other pitch detection algorithms known in the art can be used as well. On output, if the segment is  
5 determined to be unvoiced, a flag  $f_{v/uv}$  is set equal to zero and if the speech segment is voiced flag  $f_{v/uv}$  is set equal to one. Additionally, if the speech segment of delayed speech vector  $S_M(M)$  is voiced, pitch  
10 detection block 155 estimates its fundamental frequency  $F_0$  which is output to parameter encoding block 190.

In the case of an unvoiced speech segment, delayed speech vector  $S_M(M)$  is windowed in block 160  
15 by a suitable window  $w$  to generate windowed speech vector  $S_{WM}(M)$  in which the signal discontinuities to adjacent speech segments at both ends of the speech segment are reduced. Different windows, such as Hamming or Kaiser windows may be used to this end. In  
20 a specific embodiment of the present invention, a  $M$ -point normalized Hamming window  $W_H(M)$  is used, the elements of which are scaled to meet the constraint:

$$1 = \frac{1}{M} \sum_{m=0}^{M-1} w_H^2(m) \quad (1)$$

25

Windowed speech vector  $S_{WM}(M)$  is next applied to block 165 for calculating the linear prediction coding (LPC) coefficients which model the human vocal tract.  
30 As known in the art, in linear predictive coding the current signal sample  $s(n)$  is represented by a combination of the  $P$  preceding samples  $s(n-i)$ , ( $i=1, \dots, P$ ) multiplied by the LPC coefficients, plus a term which represents the prediction error. Thus, in  
35 the system of the present invention, the current

- 16 -

sample  $s(n)$  is modeled using the auto-regressive model:

$$s(n) = e_n - a_1 s(n-1) - a_2 s(n-2) - \dots - a_p s(n-p) \quad (2)$$

5

where  $a_1, \dots, a_p$  are the LPC coefficients and  $e_n$  is the prediction error. The unknown LPC coefficients which minimize the variance of the prediction error are  
 10 determined by solving a system of linear equations, as known in the art. A computationally efficient way to solve for the LPC coefficients is given by the Levinson-Durbin algorithm described for example in S.J. Orphanidis, "Optimum Signal Processing," McGraw  
 15 Hill, New York, 1988, pp. 202-207, which is hereby incorporated by reference. In a preferred embodiment of the present invention the number  $P$  of the preceding speech samples used in the prediction is set equal to 10. The LPC coefficients calculated in block 165 are  
 20 loaded into output vector  $a_{op}$ . In addition, block 165 outputs the prediction error power  $\sigma^2$  for the speech segment which is used in the decoder of the system to synthesize the unvoiced speech segment.

25 In block 170 vector  $a_{op}$ , the elements of which are the LPC coefficients, is used to solve for the roots of the homogeneous polynomial equation

$$x^n + a_1 x^{n-1} + a_2 x^{n-2} + \dots + a_{p-1} x^{n-(p-1)} + a_p = 0 \quad (3)$$

30

which roots can be recognized as the poles of the autoregressive filter modeling the human vocal tract in Eq. (2). The roots computed in block 170 are ordered in terms of increasing phase and are loaded  
 35 into pole vector  $X_p$ . The roots of the polynomial equation may be found by suitable root-finding

- 17 -

5 routines, as described for example in Press et al.,  
"Numerical Recipes, The Art of Scientific Computing,"  
Cambridge University Press, 1986, incorporated herein  
10 by reference. Alternatively, a computer implementation  
using an EISPACK set of routines can be used to  
determine the poles of the polynomial by computing the  
eigenvalues of the associated characteristic matrix,  
as used in linear systems theory and described for  
15 example in Thomas Kailath, "Linear Systems," Prentice  
Hall, Inc., Englewood Cliffs, N.J., 1980. The EISPACK  
mathematical package is described in Smith et al.,  
"Matrix Eigen System Routines - EISPACK Guide,"  
Springer-Verlag, 1976, pp. 28-29. Both publications  
are incorporated by reference.

Pole vector  $X_p$  is next received at vector  
quantizer block 180 for quantizing it into a codebook  
entry  $X_{vQ}$ . While many suitable quantization methods  
20 can be used, in a specific embodiment of the present  
invention, the quantized codebook vector  $X_{vQ}$  can be  
determined using neural networks. To this end, a  
linear vector quantizing neural network having a  
Kohonen feature map LVQ1 can be used, as described in  
25 T. Kohonen, "Self Organization and Associative  
Memory," Series in Information, Sciences, Vol. 8,  
Springer-Verlag, Berlin-Heidelberg, New York, Tokyo,  
1984, 2nd Ed. 1988.

30 It should be noted that the use of the quantized  
polynomial roots to represent the unvoiced speech  
segment is advantageous in that the dynamic range of  
the root values is smaller than the corresponding  
range for encoding the LPC coefficients thus resulting  
35 in a coding gain. Furthermore, encoding the roots of  
the prediction polynomial is advantageous in that the

- 18 -

stability of the synthesis filters can be guaranteed by restricting all poles to be less than unity in magnitude. By contrast, relatively small errors in quantizing the LPC coefficients may result in unstable poles of the synthesis filter.

The elements of the quantized  $X_{VQ}$  vector are finally input into parameter encoder 190 to form an unvoiced segment data packet for storage and transmission as described in more detail next.

In accordance with the present invention, processing of the voiced speech segments is executed in blocks 130, 140 and 150. In frame manager block 130 delayed speech vectors  $S_M(M)$  and  $S_{N-M}(N)$  are concatenated to form speech vector  $Y_N$  having a total length of  $N$  samples. In this way, an overlap of  $N-M$  samples is introduced between adjacent speech segments to provide better continuity at the segment boundaries. For voiced speech segments, the digital speech signal vector  $Y_N$  is modeled as a superposition of  $H$  harmonics expressed mathematically as follows:

$$s_N(n) = \sum_{h=0}^{H-1} A_H(h) \cdot \sin(2\pi(h+1) \frac{F_0}{f_s} n + \theta_h) + z_n; \quad (4)$$

$n=0, 1, 2, \dots, N-1.$

where  $A_H(h)$  is the amplitude corresponding to the  $h$ -th harmonic,  $\theta_h$  is the phase of the  $h$ -th harmonic,  $F_0$  and  $f_s$  are the fundamental and the sampling frequencies respectively,  $z_n$  is unvoiced noise and  $N$  is the number of samples in the enlarged speech vector  $Y_N$ .

To avoid discontinuities of the signal at the ends of the speech segments and problems associated

- 19 -

with spectral leakage during subsequent processing in the frequency domain, speech vector  $Y_N$  is multiplied in block 140 by a window  $W$  to obtain a windowed speech vector  $Y_{WN}$ . The specific window used in block 140 is a Hamming or a Kaiser window. Preferably, a  $N$  point Kaiser window  $W_K$  is used, the elements of which are normalized as shown in Eq. (1). The window functions used in the Kaiser and Hamming windows of the present invention are described in Oppenheim et al., "Discrete Time Signal Processing," Prentice Hall, Englewood Hills, NJ, 1989. The elements of vector  $Y_{WN}$  are given by the expression:

$$Y_{WN}(n) = W_K(n) \cdot Y(n); \quad n=0,1,2,\dots,N-1 \quad (5)$$

Vector  $Y_{WN}$  is received in super resolution harmonic amplitude estimation (SRHAE) block 150 which estimates the amplitudes of the harmonic frequencies on the basis of the fundamental frequency  $F_0$  of the segment obtained in pitch detector 155. The estimated amplitudes are combined into harmonic amplitude vector  $A_H$  which is input to parameter encoding block 190 to form voiced data packets.

Parameter encoding block 190 receives on input from pitch detector 155 the  $f_{vuv}$  flag which determines whether the current speech segment is voiced or unvoiced, a parameter  $E$  which is related to the energy of the segment, the quantized codebook vector  $X_{vq}$  if the segment is unvoiced, or the fundamental frequency  $F_0$  and the harmonic amplitude vector  $A_H$  if the segment is voiced. Parameter encoding block 190 outputs for each speech segment a data packet which contains all information necessary to reconstruct the speech at the receiving end of the system.

- 20 -

Figures 4 and 5 illustrate the data packets used for storage and transmission of the unvoiced and voiced speech segments in accordance with the present invention. Specifically, each data packet comprises control (synchronization) information and flag  $f_{v/uv}$  indicating whether the segment is voiced or unvoiced. In addition, each package comprises information related to the energy of the speech segment. In an unvoiced data packet this could be the sum of the squares of all speech samples or, alternatively the prediction error power computed in block 165. The information indicated as the frame energy in the voiced speech segment in Fig. 5 is preferably the sum of the estimated harmonic amplitudes computed in block 150, as described next.

As shown in Fig. 4, if the segment is unvoiced, the corresponding data packet further comprises the quantized vector  $X_{vQ}$  determined in vector quantization block 180. If the segment is voiced, the data packet comprises the fundamental frequency  $F_0$  and harmonic amplitude vector  $A_H$  from block 150, as shown in Fig. 5. The number of bits in a voiced data package is held constant and may differ from the number of bits in an unvoiced packet which is also constant.

The operation of super resolution harmonic amplitude estimation (SRHAE) block 150 is described in greater detail in Fig. 6. In step 250 the algorithm receives windowed vector  $Y_{WN}$  and the  $f_{v/uv}$  flag from pitch detector 155. In step 251 it is checked whether flag  $f_{v/uv}$  is equal to one, which indicates voiced speech. If the flag is not equal to one, in step 252 control is transferred to pole calculation block 170 (see Fig. 2). If flag  $f_{v/uv}$  is equal to one, step 253

- 21 -

is executed to determine the total number of harmonics H which is set equal to the integer number obtained by dividing the sampling frequency  $f_s$  by twice the  
 5 fundamental frequency  $F_0$ . In order to adequately represent a voiced speech segment while keeping the required bit rate low, in the system of the present invention a maximum number of harmonics  $H_{max}$  is defined and, in a specific embodiment, is set equal to 30.

10

In step 254 it is determined whether the number of harmonics H computed in step 253 is greater than or equal to the maximum number of harmonics  $H_{max}$  and if true, in step 255 the number of harmonics H is set  
 15 equal to  $H_{max}$ . In the following step 257 the input windowed vector  $Y_{WN}$  is first padded with N zeros to generate a vector  $Y_{2N}$  of length 2N defined as follows:

$$\begin{aligned} Y_{2N}(n) &= Y_{WN}(n) \quad \text{for } n=0, \dots, N-1 \\ &= 0 \quad \text{for } n=N, \dots, 2N-1 \end{aligned} \quad (6)$$

20

The zero padding operation in step 257 is required in order to obtain the discrete Fourier transform (DFT) of the windowed speech segment in vector  $Y_{WN}$  on a more finely divided set of frequencies.  
 25 It can be appreciated that dependent on the desired frequency separation, a different number of zeros may be appended to windowed speech vector  $Y_{WN}$ .

Following the zero padding, in step 257 a 2N  
 30 point discrete Fourier transform of speech vector  $Y_{2N}$  is performed to obtain the frequency domain vector  $F_{2N}$  from which the desired harmonic amplitudes are determined. Preferably, the computation of the DFT is executed using any fast Fourier transform (FFT)  
 35 algorithm of length 2N. As well known, the efficiency of the FFT computation increases if the length N of

- 22 -

the transform is a power of 2, i.e. if  $N = 2^L$ .  
 Accordingly, in a specific embodiment of the present  
 invention the length  $2N$  of the speech vector  $Y_{2N}$  may be  
 5 adjusted further by adding zeros to meet this  
 requirement. The amplitudes of the harmonic  
 frequencies of the speech segment are calculated next  
 in step 258 in accordance with the formula:

$$10 \quad A_H(h, F_0) = \frac{1}{N} \cdot \left[ 2 \cdot \sum_{k=\left[(h+1)\frac{2F_0}{f_s}N\right]-B}^{\left[(h+1)\frac{2F_0}{f_s}N\right]+B} \left[ \sum_{n=0}^{2N-1} y_{2N}(n) \cdot e^{-j2\pi\frac{k}{2N}n} \right]^2 \right]^{\frac{1}{2}}; \quad (7)$$

$$h=0, 1, 2, \dots, H-1; \quad H \leq \left\lceil \frac{f_s}{2F_0} \right\rceil$$

15 where  $A_H(h, F_0)$  is the estimated amplitude of the  $h$ -th  
 harmonic frequency,  $F_0$  is the fundamental frequency of  
 the segment and  $B$  is the half bandwidth of the main  
 lobe of the Fourier transform of the window function.

20 Considering Eq. (7) in detail we first note that  
 the expression within the inner square brackets  
 corresponds to the DFT of the windowed vector  $Y_{2N}$  which  
 is computed in step 257 and is defined as:

$$25 \quad F(k) = \sum_{n=0}^{2N-1} y_{2N}(n) e^{-j2\pi\frac{k}{2N}n} \quad (8)$$

Multiplying each resulting DFT frequency sample  
 $F(k)$  by its complex conjugate quantity  $F^*(k)$  gives the  
 power spectrum  $P(k)$  of the input signal at the given  
 30 discrete frequency sample:

$$P(k) = F(k) \cdot F^*(k) \quad (9)$$

which operation is mathematically expressed in Eq.(7)  
 by taking the square of the discrete Fourier transform  
 35 frequency samples  $F(k)$ . Finally, in Eq.(7) the  
 harmonic amplitude  $A_H(h, F_0)$  is obtained by adding



- 23 -

together the power spectrum estimates for the B adjacent discrete frequencies on each side of the respective harmonic frequency  $h$ , and taking the square  
5 root of the result, scaling it appropriately.

As indicated above, B is the half bandwidth of the discrete Fourier transform of the Kaiser window used in block 140. For a window length  $N = 512$  the  
10 main lobe of a Kaiser window has 11 samples, so that B can be rounded conveniently to 5. Since the windowing operation in block 140 corresponds in the frequency domain to the convolution of the respective transforms of the original speech segment and that of the window  
15 function, using all samples within the half bandwidth of the window transform results in an increased accuracy of the estimates for the harmonic amplitudes.

Once the harmonic amplitudes  $A_H(h, F_0)$  are  
20 computed, in step 259 the sequence of amplitudes is combined into harmonic amplitude vector  $A_H$  which is sent to the parameter encoder in step 260.

Figure 7A illustrates for comparison the harmonic  
25 amplitudes measured in an actual speech segment and the set of harmonic amplitudes estimated using the SRHAE method of the present invention. In this figure, a maximum number  $H_{max} = 30$  harmonic frequencies were used to represent an input speech segment with  
30 fundamental frequency  $F_0 = 125.36$  Hz. A normalized Kaiser window and zero padding as discussed above were also used. The percent error between the actual and estimated harmonic amplitudes is plotted in Fig. 7B and indicates very good estimation accuracy. The  
35 expression used to compute the percent error in Fig. 7B is mathematically expressed as:

- 24 -

$$E(h) = \frac{|A_a(h, F_0) - \hat{A}_a(h, F_0)|}{|A_H(h, F_0)|} \cdot 100\%; \quad \text{for } h=0, \dots, H-1. \quad (10)$$

5

The results indicate that SRHAE block 150 of the present invention is capable of providing an estimated sequence of harmonic amplitudes  $A_H(h, F_0)$  accurate to within 1000-th of a percent. Experimentally it has also been found that for a higher fundamental frequency  $F_0$ , the percent error over the total range of harmonics can be reduced even further.

10

#### B. The Decoder Block

15

Fig. 8 is a schematic block diagram of speech decoder 400 in Fig. 1. Parameter decoding block 405 receives data packets 25 via communications channel 101. As discussed above, data packets 25 correspond to either voiced or unvoiced speech segments as indicated by flag  $f_{vuv}$ . Additionally, data packets 25 comprise a parameter related to the segment energy  $E$ ; the fundamental frequency  $F_0$  and the estimated harmonic amplitudes vector  $A_H$  for voiced packets; and the quantized pole vector  $X_{VQ}$  for unvoiced speech segments.

20

25

If the current data packet 25 is unvoiced, the speech synthesis proceeds in blocks 410 through 460. Specifically, block 410 receives the quantized poles vector  $X_{VQ}$  and uses a pole codebook look up table to determine a poles vector  $X_p$  which corresponds most closely to the received vector  $X_{VQ}$ . In block 440 vector  $X_p$  is converted into a LPC coefficients vector  $a_p$  of length  $P$ . Unvoiced synthesis filter 460 is next initialized using the LPC coefficients in vector  $a_p$ . The unvoiced speech segment is synthesized by passing

30

35

- 25 -

to the synthesis filter 460 the output of white noise generator 450 which output is gain adjusted on the basis of the transmitted prediction error power  $\sigma_e$ .

5 The operation of blocks 440, 450 and 460 defining the synthesis of unvoiced speech using the corresponding LPC coefficients is known in the art and need not be discussed in further detail. Digital-to-analog converter 500 completes the process by transforming

10 the unvoiced speech segment to analog speech signal.

The synthesis of voiced speech segments and the concatenation of segments into a continuous voice signal is accomplished in the system of the present

15 invention using phase compensated harmonic synthesis block 430. The operation of synthesis block 430 is shown in greater detail in the flow diagram in Fig. 9. Specifically, in step 500 the synthesis algorithm receives input parameters from the parameter decoding

20 block 405 which includes the  $f_{v/uv}$  flag, the fundamental frequency  $F_0$  and the normalized harmonic amplitudes vector  $A_H$ . In step 510 it is determined whether the received data packet is voiced or unvoiced as indicated by the value of flag  $f_{v/uv}$ . If this value is

25 is not equal to one, in step 515 control is transferred to pole codebook search block 410 for processing of an unvoiced segment.

If flag  $f_{v/uv}$  is equal to one, indicating a voiced

30 segment, in step 520 is calculated the number of harmonics  $H$  in the segment by dividing the sampling frequency  $f_s$  of the system by twice the fundamental frequency  $F_0$  for the segment. The resulting number of harmonics  $H$  is truncated to the value of the closest

35 smaller integer.

- 26 -

Decision step 530 compares next the value of the computed number of harmonics  $H$  to the maximum number of harmonics  $H_{\max}$  used in the operation of the system.

5 If  $H$  is greater than  $H_{\max}$ , in step 540 the value of  $H$  is set equal to  $H_{\max}$ . In the following step 550 the elements of the voiced segment synthesis vector  $V_0$  are initialized to zero.

10 In step 560 the voiced/unvoiced flag  $f_{vuv}$  of previous segment is examined to determine whether the segment was voiced, in which case control is transferred in step 570 to the voiced-voiced synthesis algorithm. If the previous segment was unvoiced,  
15 control is transferred to the unvoiced-voiced synthesis algorithm. Generally, the last sample of the previous speech segment is used as the initial condition in the synthesis of the current segment as to insure amplitude continuity in the signal  
20 transition ends.

In accordance with the present invention, voiced speech segments are concatenated subject to the requirement of both amplitude and phase continuity  
25 across the segment boundary. This requirement contributes to a significantly reduced distortion and a more natural sound of the synthesized speech. Clearly, if two segments have identical number of harmonics with equal amplitudes and frequencies, the  
30 above requirement would be relatively simple to satisfy. However, in practice all three parameters can vary and thus need to be matched separately.

In the system of the present invention, if the  
35 numbers of harmonics in two adjacent voiced segments are different, the algorithm proceeds to match the

- 27 -

smallest number H of harmonics common to both segments. The remaining harmonics in any segment are considered to have zero amplitudes in the adjacent segment.

The problem of harmonics matching is illustrated in Fig. 10 where two sinusoidal signals  $s'(n)$  and  $s(n)$  having different amplitudes  $A'$  and  $A$  and fundamental frequencies  $F'_0$  and  $F_0$  have to be matched at the boundary of two adjacent segments of length  $M$ . In accordance with the present invention, the amplitude discontinuity is resolved by means of a linear amplitude interpolation such that at the beginning of the segment the amplitude of the signal  $S(n)$  is set equal to  $A'$  while at the end it is equal to the harmonic amplitude  $A$ . Mathematically this condition is expressed as

$$A'(m) + \frac{A(m) - A'(m)}{M} \quad (11)$$

where  $M$  is the length of the speech segment.

In the more general case of  $H$  harmonic frequencies the current segment speech signal may be represented as follows:

$$S(m) = \sum_{h=0}^{H-1} \left( A'(m) + \frac{A(m) - A'(m)}{M} \right) \sin((h+1)\Phi(m) + \xi(h)); \quad m=0, \dots, M-1. \quad (12)$$

where  $\Phi(m) = 2\pi m F_0/f_0$ ; and  $\xi(h)$  is the initial phase of the  $h$ -th harmonic. Assuming that the amplitudes of each two harmonic frequencies to be matched are equal, the condition for phase continuity may be expressed as an equality of the arguments of the sinusoids in Eq. (12) evaluated at the first

- 28 -

sample of the current speech segment. This condition can be expressed mathematically as:

$$\begin{aligned} (h+1) \Phi(0) + \xi(h) &= (h+1) \Phi^-(M) + \xi^-(h) \\ \xi(h) &= \Phi^-(M) + \xi^-(h); \quad \text{for } h=0, \dots, H-1 \end{aligned} \quad (13)$$

where  $\Phi$  and  $\xi$  denote the phase components for the previous segment and term  $2\pi$  has been omitted for convenience. Since at  $m = 0$  the quantity  $\Phi(m)$  is always equal to zero, Eq. (13) gives the condition to initialize the phases of all harmonics.

Fig. 11 is a flow diagram of the voiced-voiced synthesis block of the present invention which implements the above algorithm. Following the start step 600 in step 610 the system checks whether there is a DC offset  $V_0$  in the previous segment which has to be reduced to zero. If there is no such offset, in steps 620, 622 and 624 the system initializes the elements of the output speech vector to zero. If there is a DC offset, in step 612 the system determines the value of an exponential decay constant  $\gamma$  using the expression:

$$\gamma = \frac{-\log\left(\frac{0.4}{|V_0|}\right)}{M-1} \quad (14)$$

where  $V_0$  is the DC offset value.

In steps 614, 616 and 618 the constant  $\gamma$  is used to initialize the output speech vector  $S(m)$  with an exponential decay function having a time constant equal to  $\gamma$ . The elements of speech vector  $S(m)$  are given by the expression:

35

- 29 -

$$S(m) = V_0 e^{-\gamma m} \quad (15)$$

Following the initialization of the speech output  
5 vector, the system computes in steps 626, 628 and 630  
the phase line  $\phi(m)$  for time samples  $0, \dots, M$ .

In steps 640 through 670 the system synthesizes a  
segment of voiced speech of length  $M$  samples which  
10 satisfies the conditions for amplitude and phase  
continuity to the previous voiced speech segment.  
Specifically, step 640 initializes a loop for the  
computation of all  $H$  harmonic frequencies. In step  
650 the system sets up the initial conditions for the  
15 amplitude and space continuity for each harmonic  
frequency as defined in Eqs. (11)-(13) above.

In steps 660, 662 and 664 the system loops  
through all  $M$  samples of the speech segment computing  
20 the synthesized voiced segment in step 662 using  
Eq. (12) and the initial conditions set up in step  
650. When the synthesis signal is computed for all  $M$   
points of the speech segment and all  $H$  harmonic  
frequencies, following step 670 control is transferred  
25 in step 680 to initial conditions block 800.

The unvoiced-to-voiced transition in accordance  
with the present invention is determined using the  
condition that the last sample of the previous segment  
30  $S(N)$  should be equal to the first sample of the  
current speech segment  $S(N+1)$ , i.e.  $S(N) = S(N+1)$ .  
Since the current segment is voiced, it can be modeled  
as a superposition of harmonic frequencies so that the  
condition above can be expressed as:  
35 where  $A_i$  is the  $i$ -th harmonics amplitude,  $\phi_i$  and  $\theta_i$  are  
the  $i$ -th harmonics phase and initial phase,

- 30 -

$$S(N) = A_1(\phi_1 + \theta_1) + A_2(\phi_2 + \theta_2) + \dots + A_{H-1}\sin(\phi_{H-1} + \theta_{H-1}) + \xi. \quad (16)$$

5 respectively, and  $\xi$  is an offset term modeled as an exponential decay function, as described above. Neglecting for a moment the  $\xi$  term and assuming that at time  $n = N+1$  all harmonic frequencies have equal phases, the following condition can be derived:

$$10 \quad S(N) = \alpha [A_0 + A_1 + \dots + A_{H-1}] =$$

$$\alpha = \frac{S(N)}{\sum_{i=0}^{H-1} A_i} = \sin(\phi_i + \theta_i); \quad i=0, \dots, H-1. \quad (17)$$

15 where it is assumed that  $|\alpha| < 1$ . This set of equations yields the initial phases of all harmonics at sample  $n = N+1$ , which are given by the following expression:

$$\theta_i = \sin^{-1}(\alpha) - \phi_i; \quad \text{for } i=0, \dots, H-1. \quad (18)$$

20 Fig. 12 is a flow diagram of the unvoiced-voiced synthesis block which implements the above algorithm. In step 700 the algorithm starts, following an indication that the previous speech segment was unvoiced. In steps 710 to 714 the vector comprising  
25 the harmonic amplitudes of the previous segment is updated to store the harmonic amplitudes of the current voiced segment.

In step 720 a variable Sum is set equal to zero  
30 and in the following steps 730, 732 and 734 the algorithm loops through the number of harmonic frequencies H adding the estimated amplitudes until the variable Sum contains the sum of all amplitudes of the harmonic frequencies. In the following step 740,  
35 the system computes the value of the parameter  $\alpha$  after checking whether the sum of all harmonics is not equal



- 31 -

to zero. In steps 750 and 752 the value of  $\alpha$  is adjusted, if  $|\alpha| > 1$ . Next, in step 754 the algorithm computes the constant phase offset  $\beta = \sin^{-1}(\alpha)$ .

- 5 Finally, in steps 760, 762 and 764 the algorithm loops through all harmonics to determine the initial phase offset  $\theta_i$  for each harmonic frequency.

- Following the synthesis of the speech segment,
- 10 the system of the present invention stores in a memory the parameters of the synthesized segment to enable the computation of the amplitude and phase continuity parameters used in the following speech frame. The process is illustrated in a flow diagram form in Fig.
- 15 13 where in step 800 the amplitudes and phases of the harmonic frequencies of the voiced frame are loaded. In steps 810 to 814 the system updates the values of the  $H$  harmonic amplitudes actually used in the last voiced frame. In steps 820 to 824 the system sets the
- 20 values for the parameters of the unused  $H_{\max} - H$  harmonics to zero. In step 830 the voiced/unvoiced flag  $f_{v/uv}$  is set equal to one, indicating the previous frame was voiced. The algorithm exits in step 840.

- 25 The method and system of the present invention provide the capability of accurately encoding and synthesizing voiced and unvoiced speech at a minimum bit rate. The invention can be used in speech compression for representing speech without using a
- 30 library of vocal tract models to reconstruct voiced speech. The speech analysis used in the encoder of the present invention can be used in speech enhancement for enhancing and coding of speech without the use of a noise reference signal. Speech
- 35 recognition and speaker recognition systems can use the method of the present invention for modeling the

- 32 -

phonetic elements of language. Furthermore, the speech analysis and synthesis method of this invention provide natural sounding speech which can be used in artificial synthesis of a user's voice.

The method and system of the present invention may also be used to generate different sound effects. For example, changing the pitch frequency  $F_0$  and/or the harmonic amplitudes in the decoder block will have the perceptual effect of altering the voice personality in the synthesized speech with no other modifications of the system being required. Thus, in some applications while retaining comparable levels of intelligibility of the synthesized speech the decoder block of the present invention may be used to generate different voice personalities. A separate type of sound effects may be created if the decoder block uses synthesis frame sizes different from that of the encoder. In such case, the synthesized time segments will be expanded or contracted in time compared to the originals, changing their perceptual quality. The use of different frame sizes at the input and the output of an digital system, known in the art as time warping, may also be employed in accordance with the present invention to control the speed of the material presentation, or to obtain a better match between different digital processing systems.

It should further be noted that while the method and system of the present invention have been described in the context of speech processing, they are also applicable in the more general context of audio processing. Thus, the input signal of the system may include music, industrial sounds and others. In such case, dependent on the application,

- 33 -

it may be necessary to use sampling frequency higher or lower than the one used for speech, and also adjust the parameters of the filters in order to adequately represent all relevant aspects of the input signal. When applied to music, it is possible to bypass the unvoiced segment processing portions of the encoder and the decoder of the present system and merely transmit or store the harmonic amplitudes of the input signal for subsequent synthesis. Furthermore, harmonic amplitudes corresponding to different tones of a musical instrument may also be stored at the decoder of the system and used independently for music synthesis. Compared to conventional methods, music synthesis in accordance with the method of the present invention has the benefit of using significantly less memory space as well as more accurately representing the perceptual spectral content of the audio signal.

While the invention has been described with reference to a preferred embodiment, it will be appreciated by those of ordinary skill in the art that modifications can be made to the structure and form of the invention without departing from its spirit and scope which is defined in the following claims.

30

35

- 34 -

I CLAIM:

1. A method for processing an audio signal comprising the steps of:
  - 5       dividing the signal into segments, each segment representing one of a succession of time intervals;  
          detecting for each segment the presence of a fundamental frequency;  
          if such a fundamental frequency is detected,  
10       estimating the amplitudes of a set of sinusoids harmonically related to the detected fundamental frequency, the set of sinusoids being representative of the signal in the time segment; and  
          encoding for subsequent storage and transmission  
15       the set of the estimated harmonic amplitudes, each amplitude being normalized by the sum of all amplitudes.
2. The method of claim 1 wherein the audio signal  
20   is a speech signal and following the step of detecting the method further comprises the step of determining whether a segment represents voiced or unvoiced speech on the basis of the detected fundamental frequency.
3. The method of claim 2 further comprising the  
25   steps of:
  - computing a set of linear predictive coding (LPC) coefficients for each segment determined to be unvoiced; and  
30       encoding the LPC coefficients by computing the roots of a LPC coefficients polynomial.
4. The method of claim 3 further comprising the  
35   step of encoding the linear prediction error power associated with the computed LPC coefficients.

- 35 -

5. The method of claim 4 wherein the step of encoding the LPC coefficients comprises the step of computing the roots of a LPC coefficients polynomial and encoding the computed polynomial roots.

6. The method of claim 5 wherein the step of encoding the computed polynomial roots comprises the steps of: forming a vector of the computed polynomial roots; and vector quantizing the formed vector using a neural network to determine a vector codebook entry.

7. The method of claim 3 wherein each segment determined to be unvoiced is windowed with a normalized Hamming window prior to the step of computing the LPC coefficients.

8. The method of claim 2 wherein the step of estimating harmonic amplitudes comprises the steps of: performing a discrete Fourier transform (DFT) of the speech signal; and computing a root sum square of the samples of the power DFT of said speech signal in the neighborhood of each harmonic frequency to obtain an estimate of the corresponding harmonic amplitude.

9. The method of claim 8 wherein prior to the step of performing a DFT the speech signal is windowed by a window function providing reduced spectral leakage.

10. The method of claim 9 wherein the used window is a normalized Kaiser window.

- 36 -

11. The method of claim 9 wherein the computation of the DFT is accomplished using a fast Fourier transform (FFT) of the windowed segment.

5

12. The method of claim 9 wherein the estimates of the harmonic amplitudes  $A_H(h, F_0)$  are computed according to the equation:

$$10 \quad A_H(h, F_0) = \frac{1}{N} \cdot \left[ 2 \cdot \sum_{k=\left[(h+1)\frac{2F_0}{f_s}N\right]-B}^{\left[(h+1)\frac{2F_0}{f_s}N\right]+B} \left[ \sum_{n=0}^{2N-1} y_{2N}(n) \cdot e^{-j2\pi \frac{k}{2N}n} \right]^2 \right]^{\frac{1}{2}}$$

$$h=0, 1, 2, \dots, H-1; \quad H \leq \left\lceil \frac{f_s}{2F_0} \right\rceil$$

15 where  $A_H(h, F_0)$  is the estimated amplitude of the h-th harmonic frequency;  $F_0$  is the fundamental frequency; B is the half bandwidth of the main lobe of the Fourier transform of the window function; and  $y_{2N}(n)$  is the windowed input signal padded with N  
20 zeros.

13. The method of claim 12 wherein following the computation of the harmonic amplitudes  $A_H(h, F_0)$  each amplitude is normalized by the sum of all amplitudes  
25 and is encoded to obtain a harmonic amplitude vector having H elements representative of the signal segment.

14. The method of claim 5 further comprising the  
30 step of forming a data packet corresponding to each unvoiced segment for subsequent transmission or storage, the packet comprising a flag indicating that the speech segment is unvoiced, the vector codebook entry for the roots of the LPC coefficients polynomial  
35 and the linear prediction error power associated with the computed LPC coefficients.

- 37 -

15. The method of claim 13 further comprising the step of forming a data packet corresponding to each voiced segment for subsequent transmission or storage, the packet comprising a flag indicating that the speech segment is voiced, the fundamental frequency, the normalized harmonic amplitude vector and the sum of all harmonic amplitudes.

16. A method for synthesizing audio signals from data packets, at least one of the data packets representing a time segment of a signal characterized by the presence of a fundamental frequency, said at least one data packet comprising a sequence of encoded amplitudes of harmonic frequencies related to the fundamental frequency, the method comprising the steps of:

for each data packet detecting the presence of a fundamental frequency; and

synthesizing an audio signal in response only to the detected fundamental frequency and the sequence of amplitudes of harmonic frequencies in said at least one data packet.

17. The method of claim 16 wherein the audio signals being synthesized are speech signals and wherein following the step of detecting the method further comprises the steps of:

determining whether a data packet represents a voiced or unvoiced speech segment on the basis of the detected fundamental frequency;

synthesizing unvoiced speech in response to encoded information in a data packet determined to represent unvoiced speech; and

providing amplitude and phase continuity on the boundary between adjacent synthesized speech segments.

- 38 -

18. The method of claim 17 wherein the step of synthesizing unvoiced speech comprises the step of passing a white noise signal through an autoregressive digital filter the coefficients of which are the LPC coefficients corresponding to the unvoiced speech segment and the gain of the filter is adjusted on the basis of the prediction error power associated with the LPC coefficients.

10

19. The method of claim 17 wherein the step of synthesizing a voiced speech comprises the steps of:  
determining the initial phase offsets for each harmonic frequency; and

15

synthesizing voiced speech using the encoded sequence of amplitudes of harmonic frequencies and the determined phase offsets.

20

20. The method of claim 17 wherein the step of synthesizing voiced speech comprises the steps of:  
computing the frequencies of the harmonics on the basis of the fundamental frequency of the segment;  
generating voiced speech as a superposition of harmonic frequencies with amplitudes corresponding to the encoded amplitudes in the voiced data packet and phases determined as to insure phase continuity at the boundary between adjacent speech segments.

30

21. The method of claim 17 wherein the step of providing amplitude and phase continuity on the boundary between adjacent synthesized speech segments comprises the steps of:

35

determining the difference between the amplitude  $A(h)$  of  $h$ -th harmonic in the current segment and the corresponding amplitude  $A'(h)$  of the previous segment, the difference being denoted as  $\Delta A(h)$ ; and



- 39 -

providing a linear interpolation of the current segment amplitude between the end points of the segment using the formula:

$$5 \quad A(h,m) = A(h,0) + m \cdot \Delta A(h) / M, \quad \text{for } m = 0, \dots, M-1.$$

22. The method of claim 19 wherein the voiced speech is synthesized using the equation:

$$10 \quad S(m) = \sum_{h=0}^{H-1} \left( A(h) + \frac{\Delta A(h)}{M} \cdot m \right) \sin((h+1)\phi(m) + \xi(h));$$

$$m=0, \dots, M-1. \quad (20)$$

where  $A(h)$  is the amplitude of the signal at the end of the previous segment;  $\phi(m) = 2\pi m F_0 / f_s$ , where  $F_0$  is the fundamental frequency and  $f_s$  is the sampling frequency; and  $\xi(h)$  is the initial phase of the  $h$ -th harmonic.

23. The method of claim 22 wherein phase continuity for each harmonic frequency in adjacent voiced segments is insured using the boundary condition:

$$\xi(h) = (h+1)\phi(M) + \xi(h),$$

where  $\phi(M)$  and  $\xi(h)$  are the corresponding quantities of the previous segment.

24. The method of claim 22 wherein the initial phase for each harmonic frequency in an unvoiced-to-voiced transition is computed using the condition:

$$30 \quad \xi(h) = \sin^{-1}(\alpha);$$

$$\alpha = \frac{S(M)}{\sum_{i=0}^{H-1} A_i}; \quad i=0, \dots, H-1.$$

35

- 40 -

where  $S(M)$  is the  $M$ -th sample of the unvoiced speech segment;  $A_i$  are the harmonic amplitudes for  $i = 0, \dots, H-1$ ; and  $|\alpha| < 1$ , and  $\phi(m)$  is evaluated at the  $M+1$  sample.

25. The method of claim 24 further comprising the step of generating sound effects by changing the fundamental frequency  $F_0$  and the values of the harmonic amplitudes encoded in the data packet.

26. The method of claim 24 further comprising the step of generating sound effects by changing the length of the synthesized signal segments.

15

27. A system for processing audio signals comprising:

means for dividing an audio signal into segments, each segment representing one of a succession of time intervals;

20

means for detecting for each segment the presence of a fundamental frequency;

means for estimating the amplitudes of a set of sinusoids harmonically related to the detected fundamental frequency, the set of sinusoids being representative of the signal in the time segment; and

25

means for encoding the set of harmonic amplitudes, each amplitude being normalized by the sum of all amplitudes.

30

28. The system of claim 27 wherein the audio signal is a speech signal and the system further comprises means for determining whether a segment represents voiced or unvoiced speech on the basis of the detected fundamental frequency.

35

- 41 -

29. The system of claim 28 further comprising:  
means for computing a set of linear predictive  
coding (LPC) coefficients corresponding to a speech  
5 segment; and  
means for encoding the LPC coefficients and the  
linear prediction error power associated with the  
computed LPC coefficients.

10 30. The system of claim 29 wherein the means for  
encoding the LPC coefficients comprises means for  
computing the roots of a LPC coefficients polynomial  
and means for encoding polynomial roots into a  
codebook entry.

15 31. The system of claim 30 wherein the means for  
encoding polynomial roots comprises a neural network  
providing the capability of vector quantizing the  
polynomial roots into a vector codebook entry.

20 32. The system of claim 28 further comprising  
windowing means providing the capability of  
multiplying the signal segment with the coefficients  
of a predetermined window function.

25 33. The system of claim 28 wherein the means for  
estimating harmonic amplitudes comprises:  
means for performing a discrete Fourier transform  
(DFT) of a digitized signal segment; and  
30 means for computing a root sum square of the  
samples of the DFT in the neighborhood of a harmonic  
frequency, said means obtaining an estimate of the  
amplitude of the harmonic frequency.

35 34. The system of claim 33 wherein the means for  
performing a DFT computation comprises means for

- 42 -

performing a fast Fourier transform (FFT) of the signal segment.

5        35. The system of claim 33 further comprising means for padding the input signal with zeros.

10       36. The system of claim 33 further comprising means for normalizing the computed harmonic amplitudes.

15       37. The system of claim 36 further comprising means for forming a data packet corresponding to each unvoiced segment, the packet comprising a flag indicating that the speech segment is unvoiced, the codebook entry for the roots of the LPC coefficients polynomial and the linear prediction error power associated with the computed LPC coefficients; and means for forming a data packet corresponding to each voiced segment for subsequent transmission or storage, the packet comprising a flag indicating that the speech segment is voiced, the fundamental frequency, a vector of the normalized harmonic amplitudes and the sum of all harmonic amplitudes.

25       38. A system for synthesizing audio signals from data packets, at least one of the data packets representing a time segment of a signal characterized by the presence of a fundamental frequency, said at least one data packet comprising a sequence of encoded amplitudes of harmonic frequencies related to the fundamental frequency, the system comprising:

30       means for determining the fundamental frequency of the signal represented by said at least one data packet;

35

- 43 -

means for synthesizing an audio signal segment in response to the determined fundamental frequency and the sequence of amplitudes of harmonic frequencies in said at least one data packet; and

means for providing amplitude and phase continuity on the boundary between adjacent synthesized audio signal segments.

39. The system of claim 38 wherein the means for synthesizing comprises means for determining the initial phase offsets for each harmonic frequency.

40. The system of claim 39 wherein the means for providing amplitude and phase continuity comprises means for providing a linear interpolation between the values of the amplitude of the signal at the end points of the segment.

41. The system of claim 39 wherein the means for providing amplitude and phase continuity further comprises means for computing conditions for phase continuity between harmonic frequencies in adjacent speech segments in accordance with the formula:

$$\xi(h) = (h+1)\phi(M) + \xi'(h),$$

where  $\xi(h)$  is the initial phase of the h-th harmonic of the current segment;  $\phi(m) = 2\pi m F_0/f_s$ , where  $F_0$  is the fundamental frequency and  $f_s$  is the sampling frequency; and  $\xi'(M)$  and  $\xi'(h)$  are the corresponding quantities of the previous segment.

42. The system of claim 41 further comprising means for generating sound effects by changing the fundamental frequency  $F_0$  and the encoded values of the harmonic amplitudes.

- 44 -

43. The system of claim 41 further comprising means for generating sound effects by changing the size of synthesized signal segments.

5

44. A system for synthesizing speech from data packets, the data packets representing voiced or unvoiced speech segments, comprising:

- 10 means for determining whether a data packet represents a voiced or unvoiced speech segment;
- means for synthesizing unvoiced speech in response to encoded information in an unvoiced data packet;
- 15 means for synthesizing voiced speech segment signal in response only to a sequence of amplitudes of harmonic frequencies encoded in a voiced data packet;
- and
- means for providing amplitude and phase continuity on the boundary between adjacent
- 20 synthesized speech segments.

45. The system of claim 44 wherein the means for synthesizing unvoiced speech comprises: means for generating white noise; a digital synthesis filter;

25 means for initializing the coefficients of the synthesis filter using a set of parameters representative of an unvoiced speech segment, and means for adjusting the gain of the synthesis filter.

30 46. The system of claim 44 wherein the means for synthesizing a voiced speech segment comprises means for determining the initial phase offsets for each harmonic frequency.

35

- 45 -

47. The system of claim 44 wherein the means for  
providing amplitude and phase continuity comprises  
means for providing a linear interpolation between the  
5 values of the signal amplitude at the end points of  
the segment.

10

15

20

25

30

35

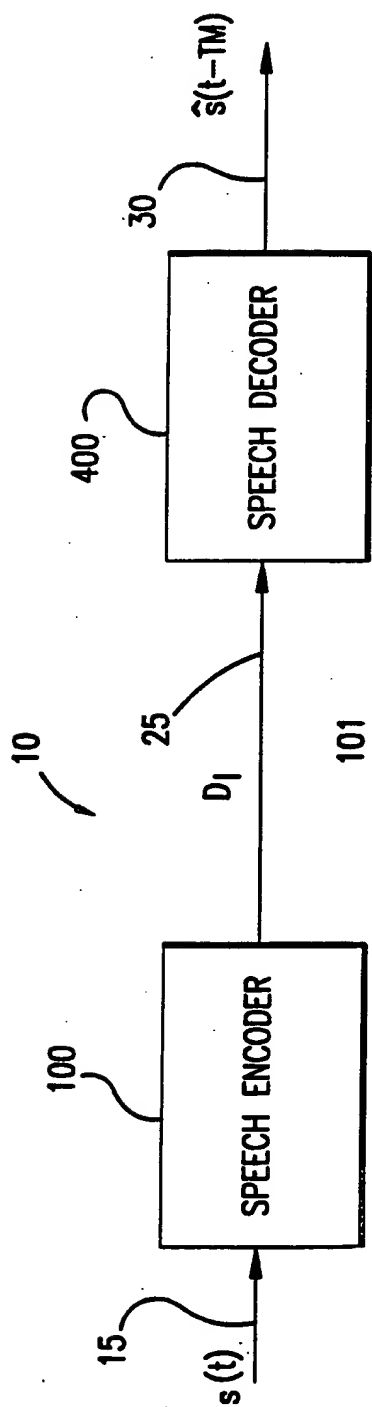


FIG. 1

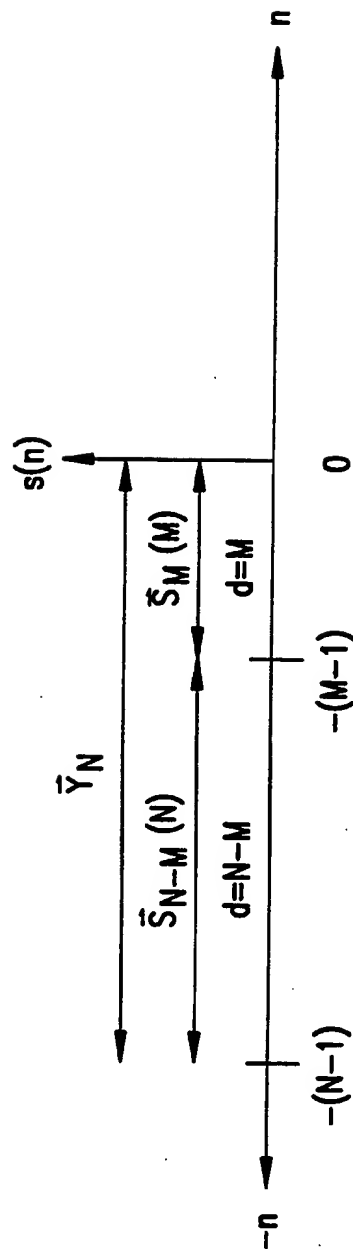


FIG. 3



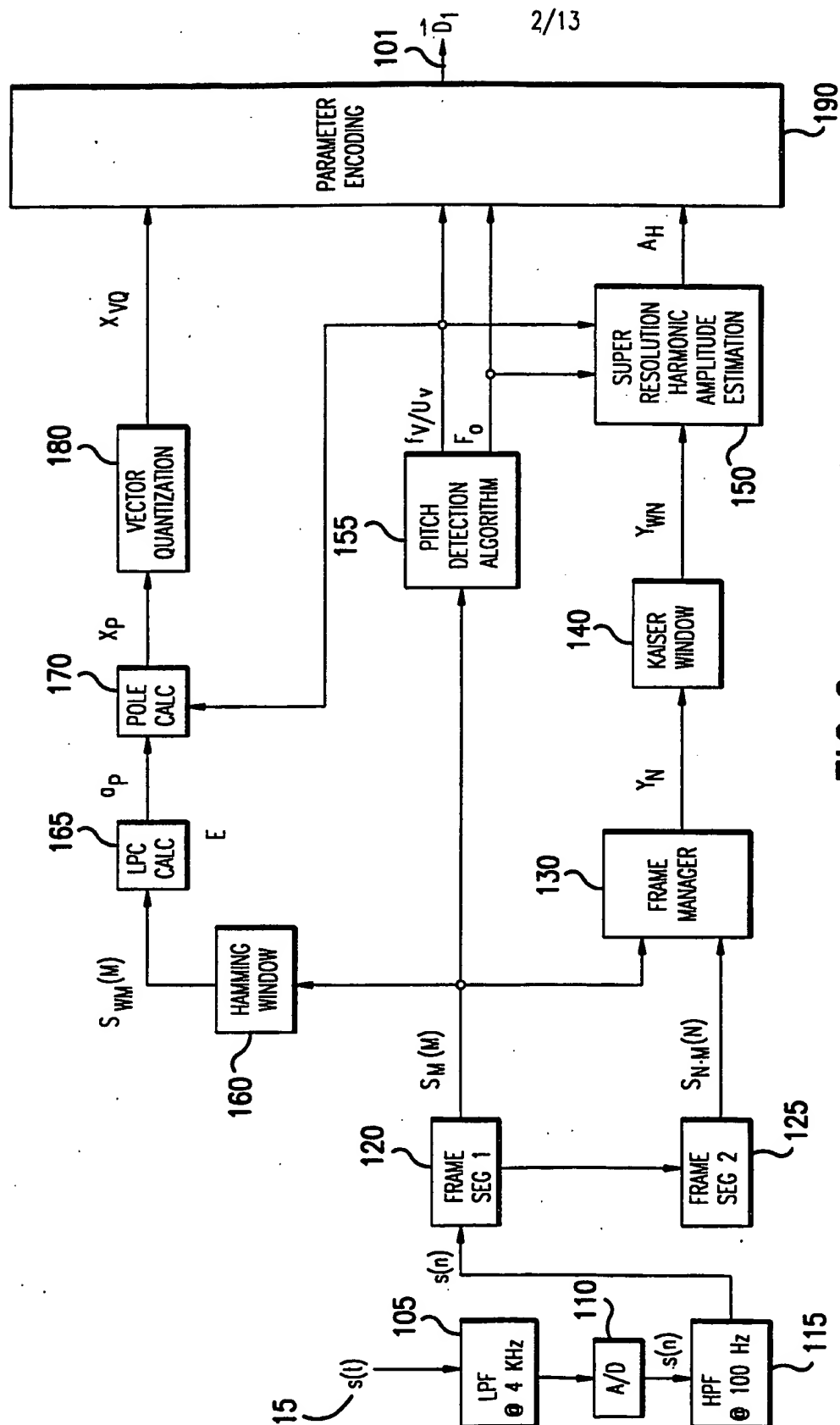


FIG. 2

SYNC	$f_v/uv$	PREDICTION POWER ERROR	CODEWORD VECTOR $X_{VQ}$
------	----------	------------------------	--------------------------

FIG.4

SYNC	$f_v/uv$	FUNDAMENTAL FREQUENCY $F_0$	FRAME ENERGY $E$	NORMALIZED HARMONIC AMPLITUDES VECTOR
------	----------	--------------------------------	---------------------	---------------------------------------

FIG.5

4/13

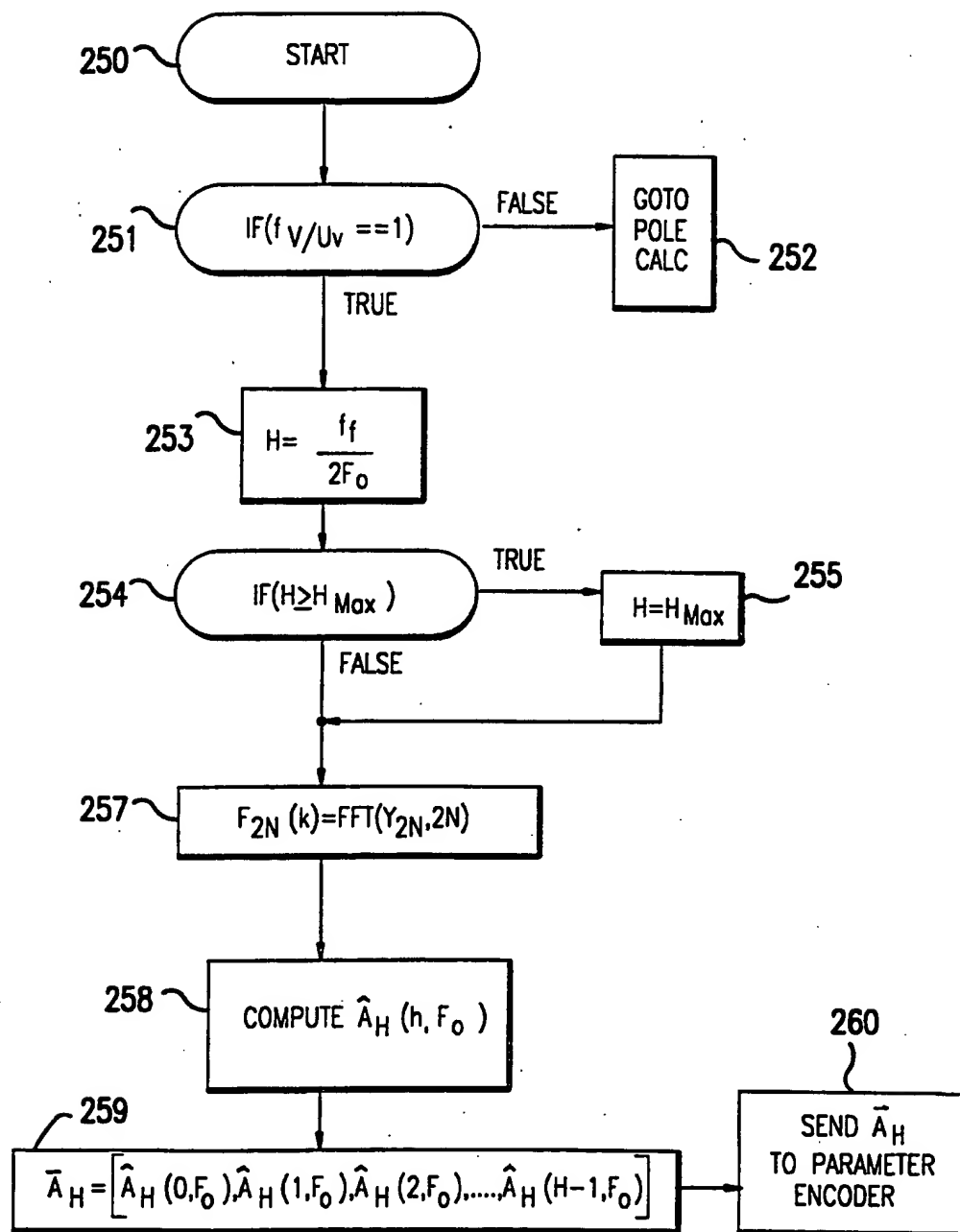


FIG. 6

5/13

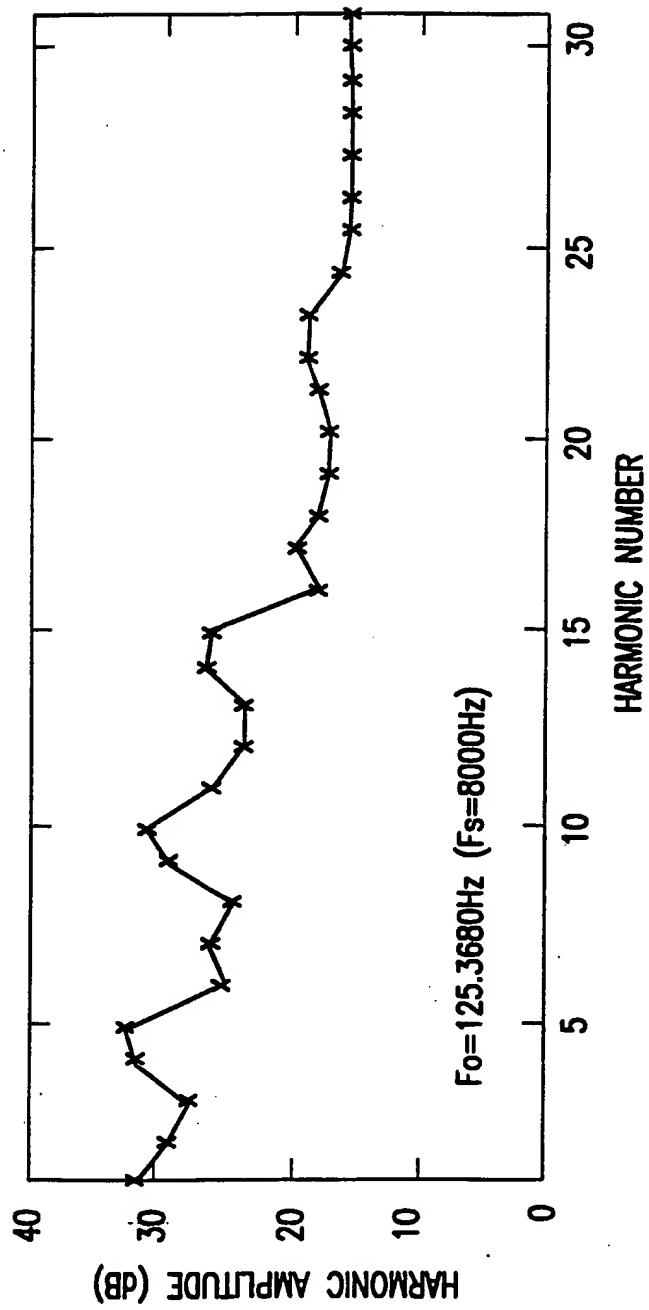


FIG.7A

6/13

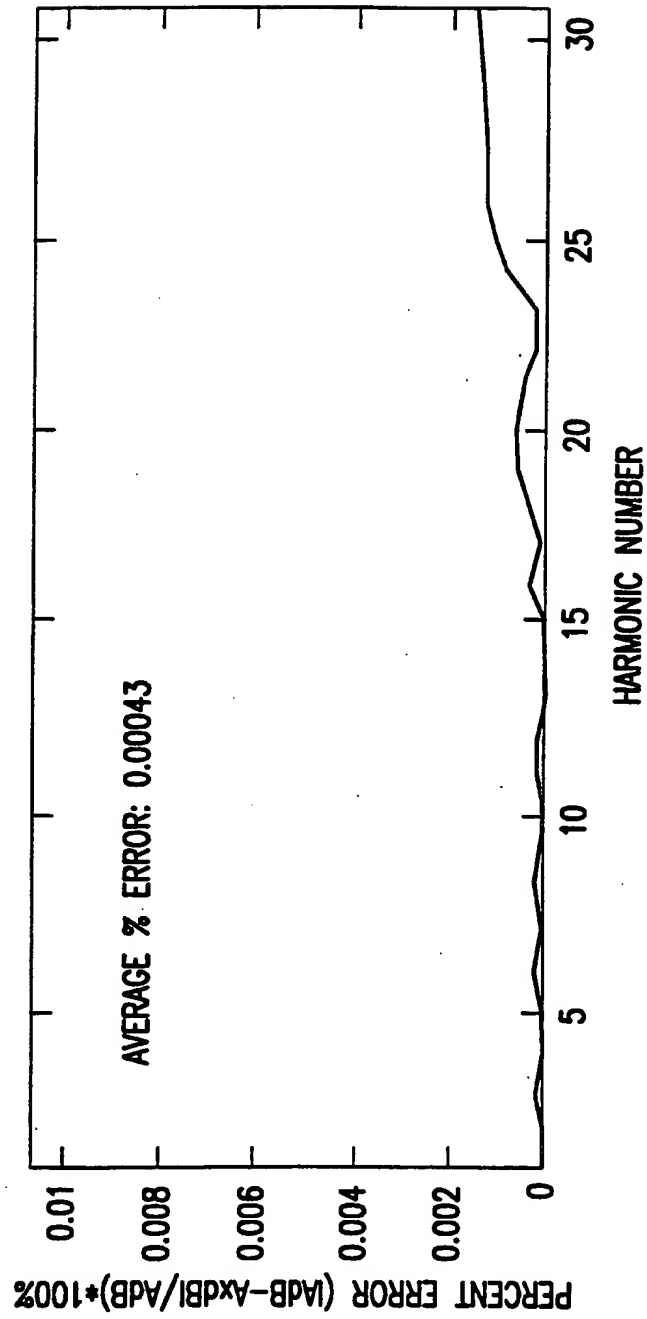
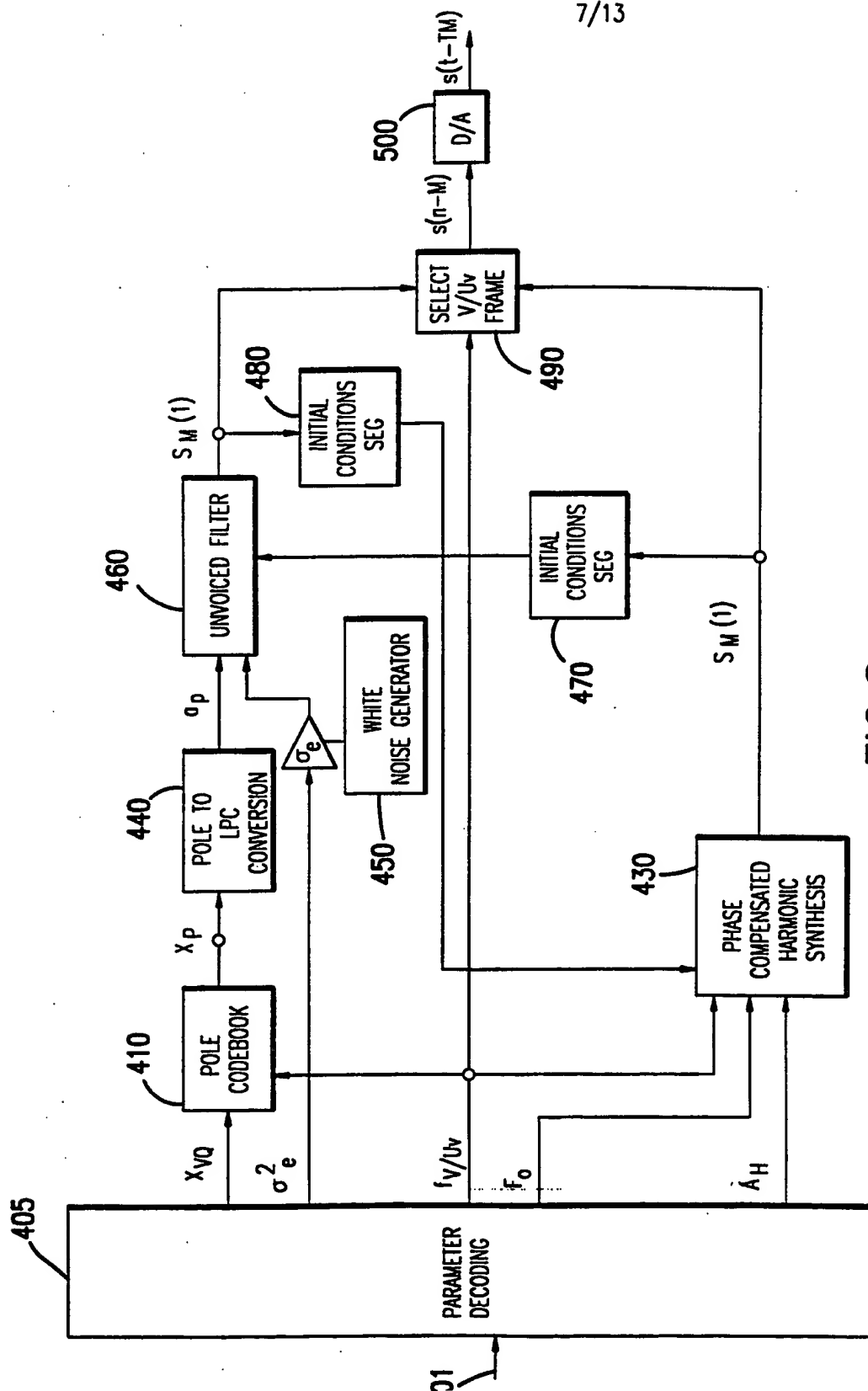


FIG.7B

7/13



8/13

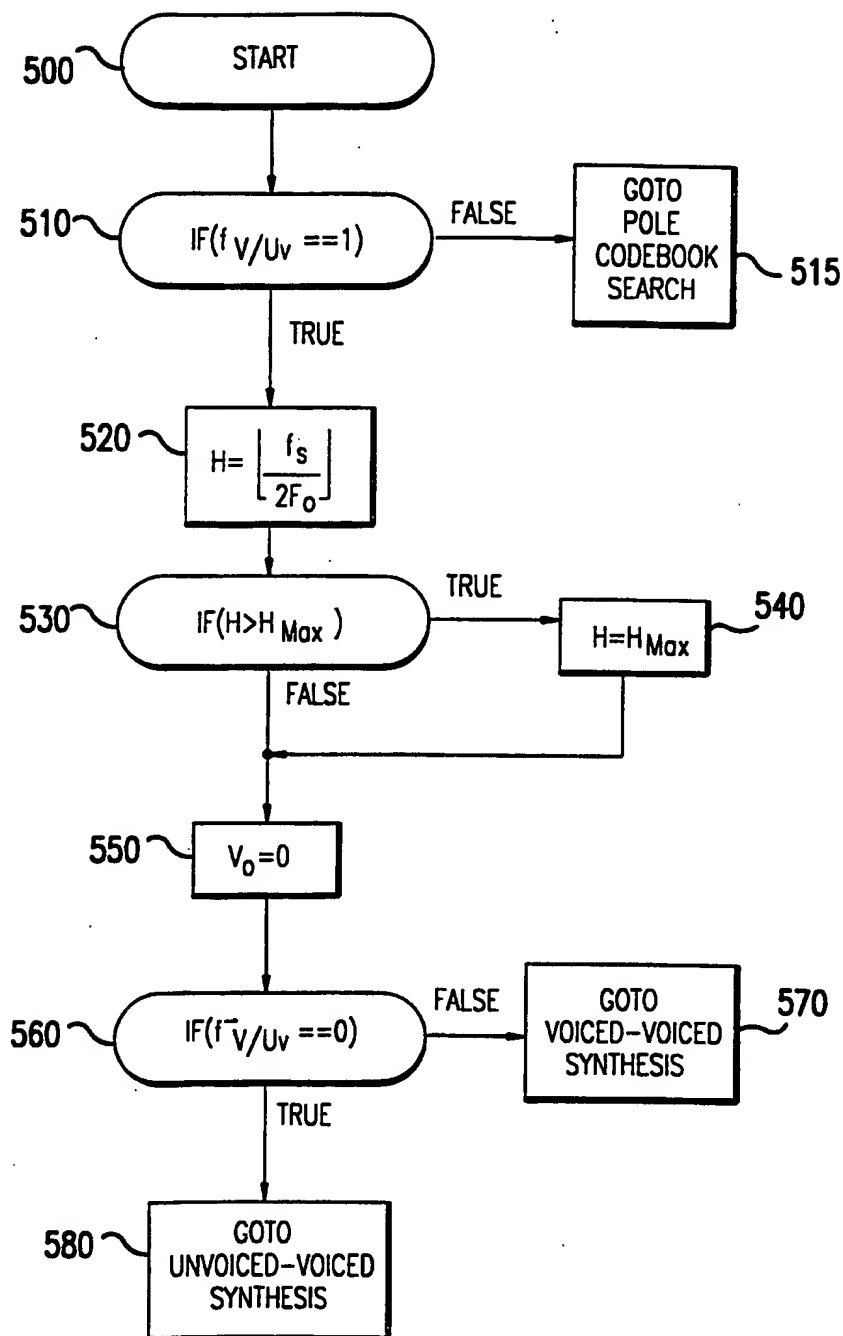


FIG. 9

9/13

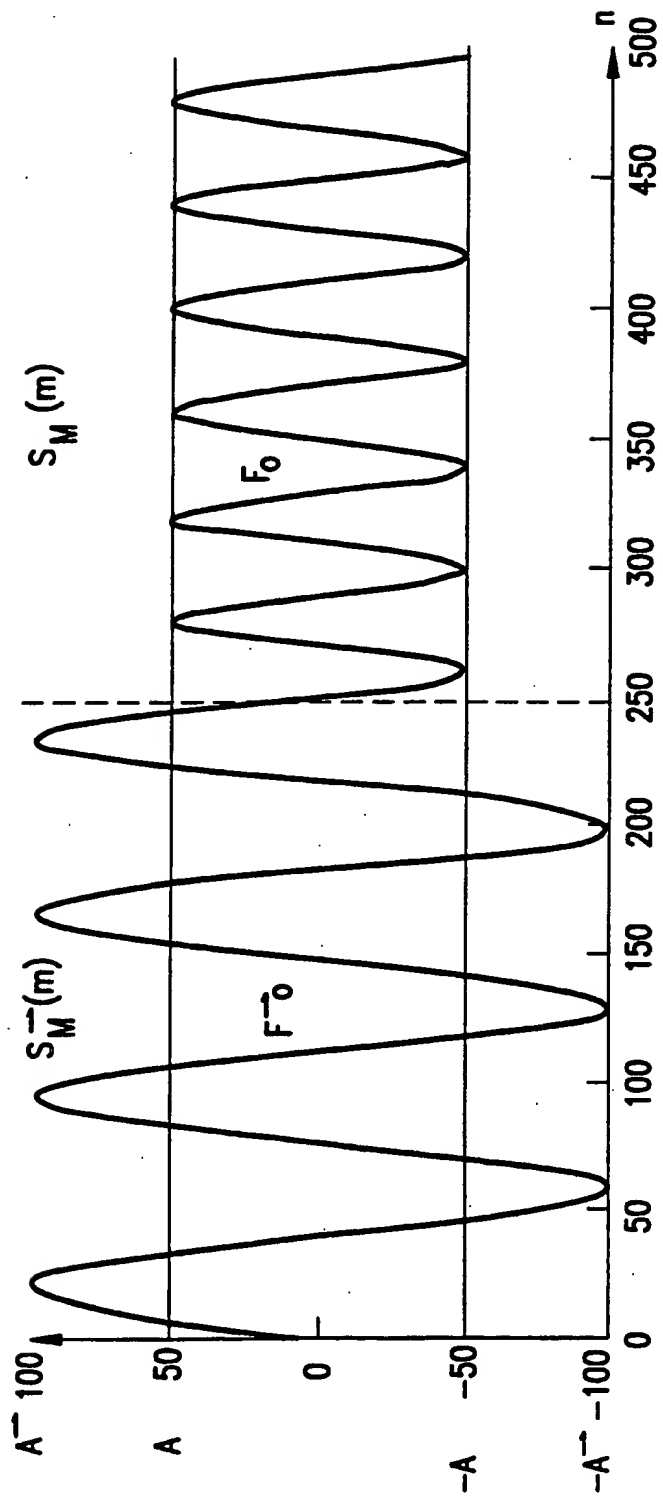


FIG.10a



10/13

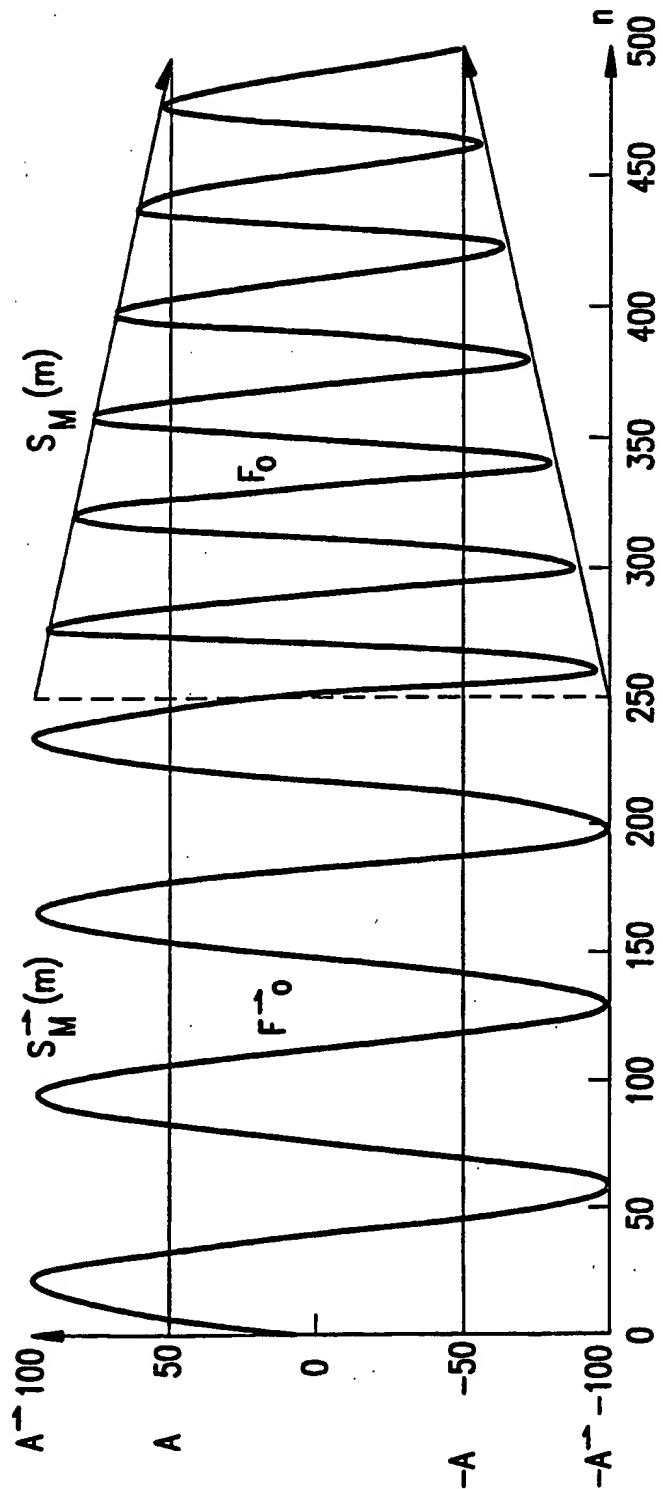


FIG.10b

11 / 13

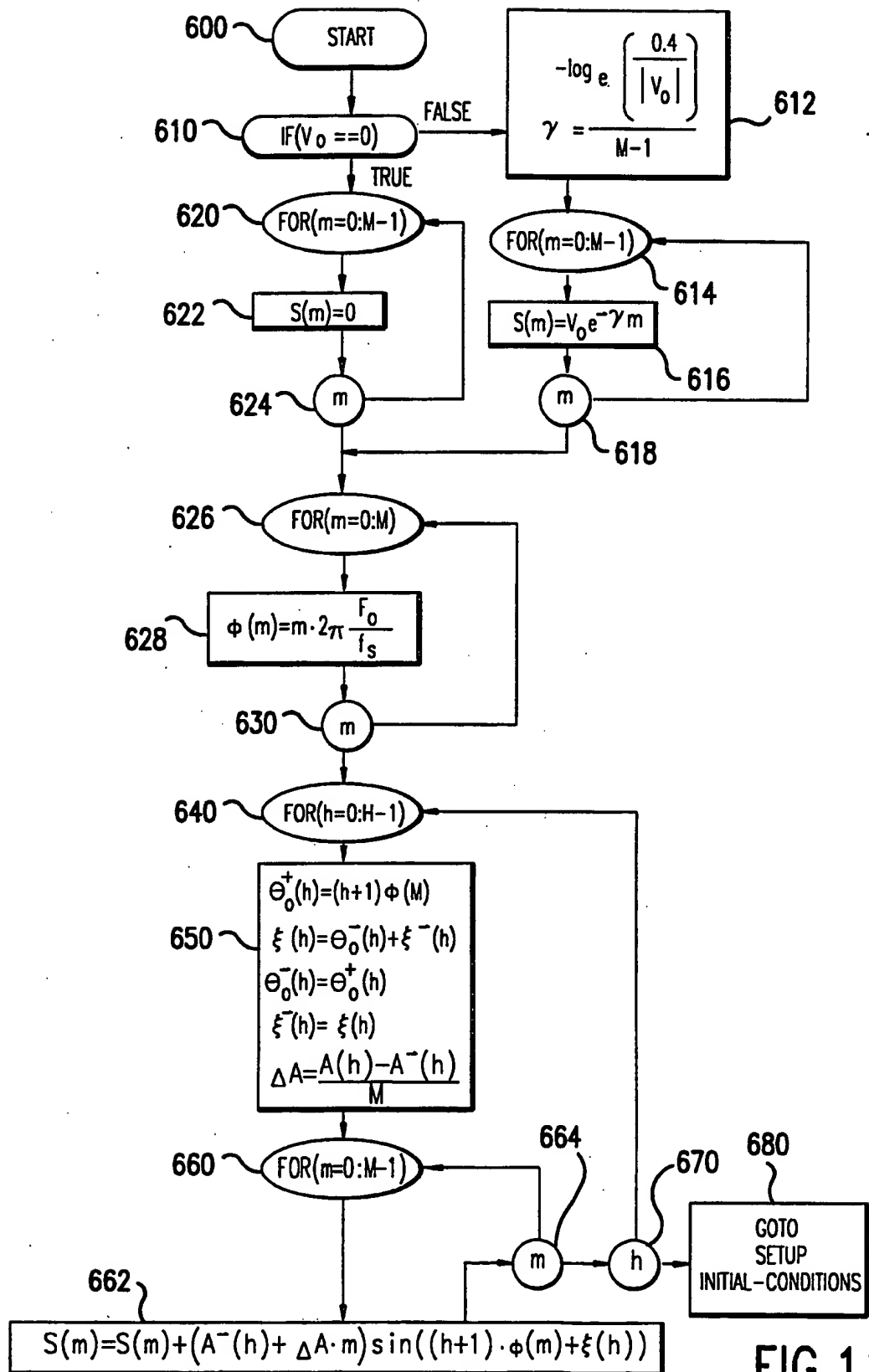


FIG.11

12 / 13

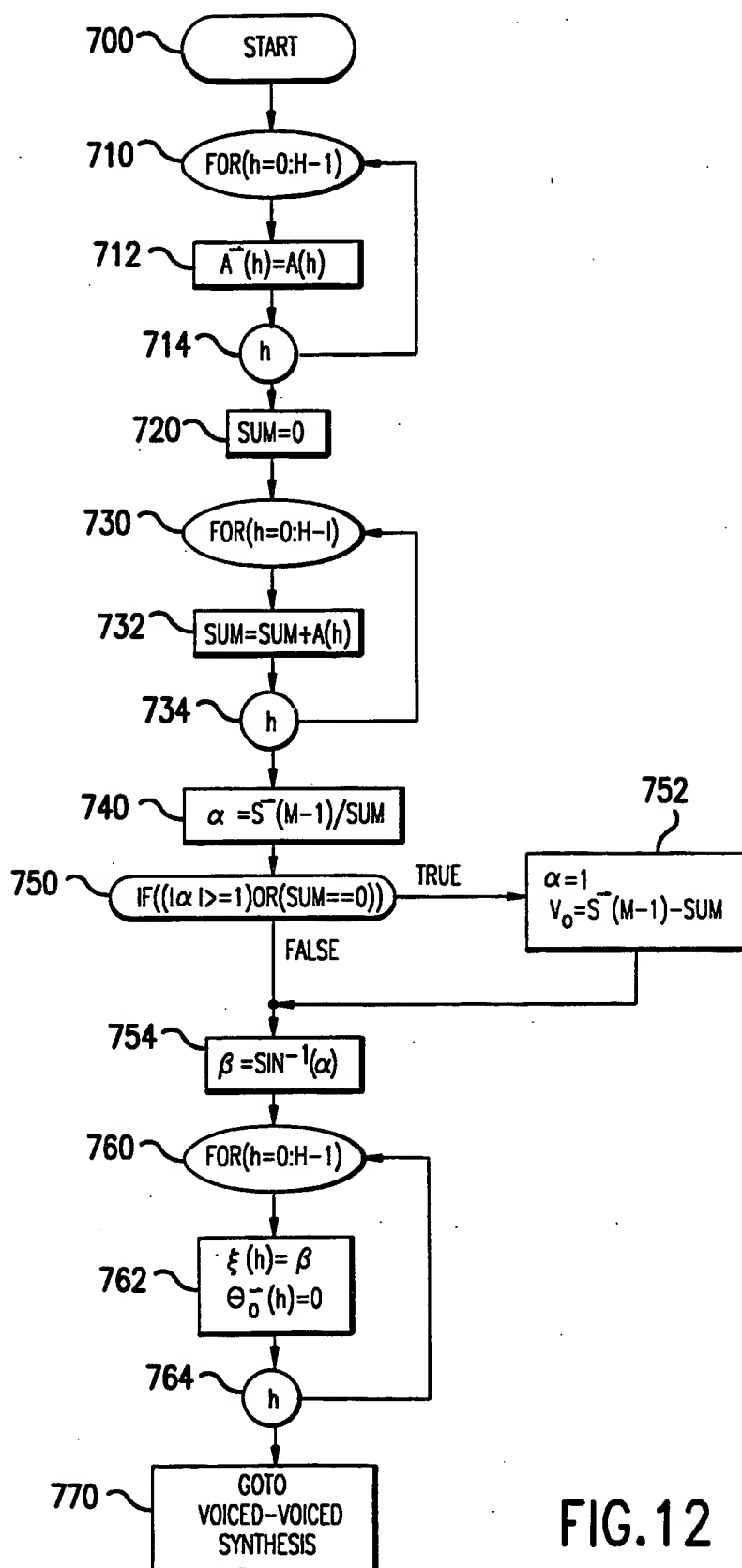


FIG. 12

13 / 13

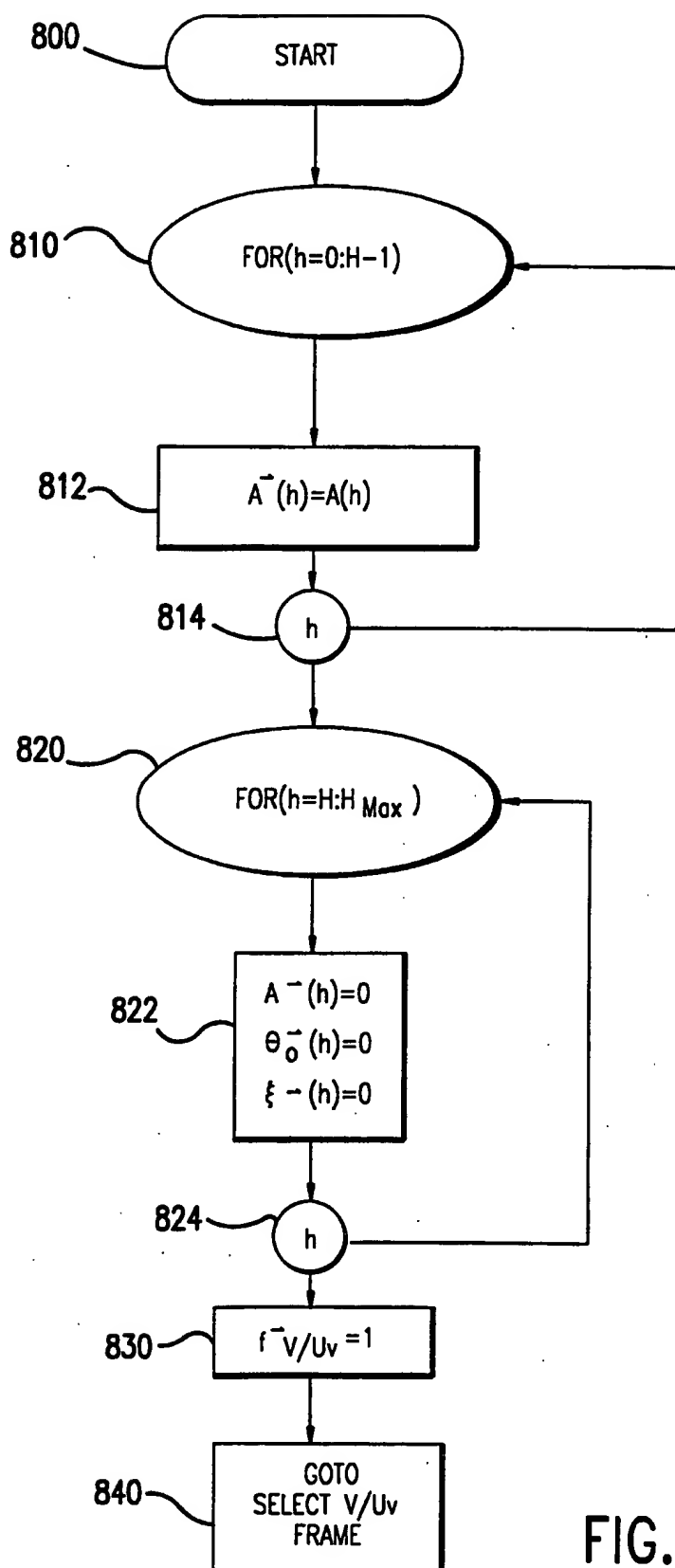


FIG.13

## INTERNATIONAL SEARCH REPORT

International application No.  
PCT/US95/08616

<b>A. CLASSIFICATION OF SUBJECT MATTER</b> IPC(6) :G10L 3/02, 9/00 US CL :395/2.17 According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b> Minimum documentation searched (classification system followed by classification symbols) U.S. : 395/2.17, 2.14, 2.15, 2.28, 2.31, 2.33, 2.71, 2.74, 2.77  Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched  Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) APS, IEEE CDROM Library search terms:voice, unvoice, segments, lpc, interpolation, Hamming window, Kaiser window, vector quantization		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US, A, 4,797,926 (BRONSON ET AL) 10 January 1989, see fig. 1.	1-47
X	US, A, 4,771,465 (BRONSON ET AL.) 13 September 1988, see fig. 1 and fig. 2.	1-47
Y	US, A, 4,435,832 (ASADA ET AL.) 06 March 1984, see abstract.	16-26, 38-47
X	US, A, 4,802,221 (JIBBE) 31 January 1989, see abstract, fig. 3.	1-15, 27-37
A	US, A, 4,864,620 (BIALICK) 05 September 1989, see abstract.	1-47
X	US, A, 4,991,213 (WILSON) 05 February 1991, see fig. 3.	1-15, 27-37
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
* Special categories of cited documents:	"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention	
"A" document defining the general state of the art which is not considered to be part of particular relevance	"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone	
"E" earlier document published on or after the international filing date	"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art	
"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	"G" document member of the same patent family	
"O" document referring to an oral disclosure, use, exhibition or other means		
"P" document published prior to the international filing date but later than the priority date claimed		
Date of the actual completion of the international search 01 SEPTEMBER 1995		Date of mailing of the international search report 15 SEP 1995
Name and mailing address of the ISA/US Commissioner of Patents and Trademarks Box PCT Washington, D.C. 20231 Facsimile No. (703) 305-3230		Authorized officer RICHMOND DORVIL Telephone No. (703) 305-9645

## INTERNATIONAL SEARCH REPORT

International application No.  
PCT/US95/08616

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US, A, 5,189,701 (JAIN) 23 February 1993, see fig. 1, fig. 3, and abstract.	1-47
A	US, A, 5,247,579 (HARDWICK ET AL) 21 September 1993, see abstract.	1-47
X	US, A, 5,303,346 (FESSELER ET AL) 12 April 1994, see abstract.	1-15, 27-37
Y	US, A, 5,327,521 (SAVIC ET AL.) 05 July 1994, see fig. 2.	1-47
X, P	US, A, 5,369,724 (LIM) 29 November 1994, see fig. 4A.	1-47
A	IEEE, Proceedings of ICASSP 1986, Tokyo, McAulay et al., "Phase Modeling and its Application Sinusoidal Transform Coding", pp.370-373.	1-47
X	IEEE, Proceedings of ICASSP 1988, Thompson, "Parametric Models of the Magnitude/Phase Spectrum for Harmonic Speech Coding", pp. 378-381.	1-47